

# **Title: Candidate gene scan for Single Nucleotide Polymorphisms involved in the determination of normal variability in human craniofacial morphology**

**Authors:** Mark Barash<sup>1,2\*</sup>, Philipp E. Bayer<sup>3,4</sup>, Angela van Daal<sup>2</sup>

## **Affiliations:**

<sup>1</sup>School of Mathematical and Physical Sciences, Centre for Forensic Sciences, Faculty of Science, University of Technology Sydney, Sydney, NSW, Australia

<sup>2</sup>Faculty of Health Sciences and Medicine, Bond University, Gold Coast, QLD, Australia

<sup>3</sup>School of Agriculture and Food Sciences, The University of Queensland, Brisbane, QLD, Australia

<sup>4</sup>Australian Centre for Plant Functional Genomics, The University of Queensland, Brisbane, QLD, Australia

Emails: Mark Barash [mark.barash@uts.edu.au](mailto:mark.barash@uts.edu.au); Philipp E. Bayer [philipp.bayer@uqconnect.edu.au](mailto:philipp.bayer@uqconnect.edu.au); Angela van Daal [vandaalconsult@gmail.com](mailto:vandaalconsult@gmail.com)

- \* Corresponding author email: [mark.barash@uts.edu.au](mailto:mark.barash@uts.edu.au)
- \* Corresponding author postal address: Centre for Forensic Sciences, School of Mathematical and Physical Sciences, Faculty of Science University of Technology Sydney, Building 4, Thomas St, Broadway NSW 2007, Australia.

# Abstract

## Background

Despite intensive research on genetics of the craniofacial morphology using animal models and human craniofacial syndromes, the genetic variation that underpins normal human facial appearance is still largely unknown. Recent development of novel digital methods for capturing the complexity of craniofacial morphology in conjunction with high-throughput genotyping methods, show great promise for unravelling the genetic basis of such a complex trait. Better understanding of the craniofacial genetics would allow development of novel tools for medical diagnostics of craniofacial syndromes as well as prediction of the facial appearance for forensic and intelligence use.

## Results

We selected 1,319 candidate craniofacial genetic markers and previously described pigmentation polymorphisms as well as additional 4,732 markers in linkage disequilibrium (LD) with candidate markers, which were subsequently genotyped using massively parallel sequencing. We manually allocated 32 craniofacial landmarks and calculated 92 craniofacial distances from 3-Dimensional (3D) facial scans and made six direct cranial measurements of 587 volunteers. We also recorded information on two facial traits (eyelid and ear lobe) and calculated ten principal components, based on all the craniofacial measurements. Genetic association between 104 craniofacial phenotypes and 3,073 genetic markers were tested. Following application of a genomic wide association study (GWAS) p-value threshold of  $5.00E-08$ , 45 single nucleotide polymorphisms (SNPs) in 27 genes and 13 intergenic regions were associated with 11 craniofacial traits. Following subsequent application of an over-conservative Bonferroni correction, associations were observed between 8 craniofacial traits and 12 SNPs located in 12 genes and intergenic regions. We report all the significant associations that reached the  $5.00E-08$  p-value threshold as we believe this threshold is conservative enough to avoid or at least significantly reduce potentially spurious associations. Majority of associations were in novel genes, while one SNP was in the PAX3 gene that was previously linked to normal variation in craniofacial morphology. Another two polymorphisms were found in the COL11A1 gene, which was previously linked to normal variation in craniofacial morphology, including confirmation of one of the associated SNPs.

Associations of the pigmentation traits in this study have fully confirmed previously published results.

## Conclusions

This study identified the greatest number of genetic variants associated with normal variation of craniofacial morphology to date by using a candidate gene approach. These results enhance our understanding of the genetics that determines normal variation in craniofacial morphology and will also be of particular value in the forensic field to allow prediction of a person's appearance from a DNA sample.

## Keywords

SNPs, single nucleotide polymorphisms, craniofacial, facial appearance, embryogenetics forensic DNA phenotyping, facial reconstruction.

## Background

The human face is probably the most commonly used descriptor of a person and has an extraordinary role in human evolution, social interactions, clinical applications as well as forensic investigations. The influence of genes on facial appearance can be seen in the striking resemblance of monozygotic twins as well as amongst first degree relatives, indicating a high heritability [1, 2].

Uncovering the genetic background for regulation of craniofacial morphology is not a trivial task. Human craniofacial development is a complex multistep process, involving numerous signalling cascades of factors that control neural crest development, followed by a number of epithelial-mesenchymal interactions that control outgrowth, patterning and skeletal differentiation, as reviewed by Sperber et. al. [2]. The mechanisms involved in this process include various gene expression and protein translation patterns, which regulate cell migration, positioning and selective apoptosis, subsequently leading to development of specific facial prominences. These events are precisely timed and are under hormonal and metabolic control. Most facial features of the human embryo are recognizable from as early as 6 weeks post conception, developing rapidly *in utero* and continuing to develop during childhood and adolescence [3, 4]. Development of the face and brain are interconnected and

occur at the same time as limb formation. Facial malformations therefore, frequently occur with brain and limb abnormalities and vice versa. Genetic regulation of craniofacial development involves several key morphogenic factors such as HOX, WNT, BMP, FGF as well as hundreds of other genes and intergenic regulatory regions, incorporating numerous polymorphisms [2]. The SNPs involved in craniofacial diseases may in fact influence the extraordinary variety of human facial appearances, in the same way that genes responsible for albinism have been shown to be involved in normal pigmentation phenotypes [5]. Additionally, non-genetic components such as nutrition, climate and socio-economic environment may also affect human facial morphology via epigenetic regulation of transcription, translation and other cellular mechanics. To date, both the genetic and even more so, the epigenetic regulation of craniofacial morphology shaping are poorly understood.

The genetic basis of craniofacial morphogenesis has been explored in numerous animal models with multiple loci shown to be involved [2]. The majority of human studies in this field have focused on the genetics of various craniofacial disorders such as craniosynostosis and cleft lip/palate [6, 7], which may provide a link to regulation of normal variation of the craniofacial phenotype, as for example observed between cleft-affected offspring and the increase of facial width seen in non-affected parents [8]. These studies have identified several genes with numerous genetic variants that may contribute to normal variation of different facial features, such as cephalic index, bizygomatic distance and nasal area measurements [9-11]. Studies of other congenital disorders involving manifestation of craniofacial abnormalities such as Alagille syndrome (JAG1 and NOTCH2 gene mutations), Down syndrome (chromosome 21 trisomy - multiple genes), Floating-Harbor syndrome (SRCAP gene mutations) and Noonan syndrome (mutations in various genes such as PTPN11 and RAF1) provide additional information on the candidate genes potentially involved in normal craniofacial development [12-17].

In recent years, new digital technologies such as 3-Dimensional laser imaging have been used in numerous anthropometric studies. 3-D laser imaging allows accurate and rapid capture of facial morphology, providing a better alternative to traditional manual measurements of craniofacial distances [18-20]. The high-throughput genotyping technologies and digital methods for capturing facial morphology have been used in a number of recent studies that demonstrated a link between normal facial variation and specific genetic polymorphisms [21-23]. Despite these promising results, our current knowledge of craniofacial genetics is sparse.

This study aims to further define the polymorphisms associated with normal facial variation using a candidate gene approach. The advantage of a candidate gene approach over previous genome wide association studies (GWAS) is that it focuses on genes, which have previously been associated with craniofacial embryogenesis or inherited craniofacial syndromes, rather than screening hundreds of thousands of non-specific markers. This approach aims to increase the chances of finding significant associations between SNPs and visible traits and requires fewer samples for robust association analysis [24, 25].

In the current study, 32 anthropometric landmarks and 92 craniofacial trait measurements were made from the 3-D facial scans of 587 volunteers of various ancestries. The calculation of principal components based on the craniofacial measurements was performed in order to obtain a more comprehensive representation of the facial shape. The associations between these craniofacial traits and 3,073 genetic markers were tested.

This research should assist in uncovering the genetic basis of normal craniofacial morphology variation and will enhance our understanding of craniofacial embryogenetics. These findings could be useful in building mathematical models to predict facial appearance from a forensic DNA sample where no suspect has been identified, thereby providing valuable investigative leads. It could also assist in identifying skeletal remains by allowing more accurate facial reconstructions.

# Methods

## Sample collection and ethics statement

A total of 623 unrelated individuals, mostly Bond University (Gold Coast, Australia) students, of diverse ancestry backgrounds were recruited. The participants provided their written informed consent to participate in this study, which was approved by the Bond University Ethics committee (RO-510). To minimize any age-related influences on facial morphology the samples were largely collected from volunteers aged between 18 and 40. The mean age of the volunteers was 26.6 (SD  $\pm$  8.9). Individuals who had experienced severe facial injury and/or undergone facial surgery (e.g. nose or chin plastics) were not included in analysis of the facial traits.

Each participant donated four buccal swabs (Isohelix, Cell Projects, Kent, UK). 3-Dimensional (3-D) facial scans and three direct cranial measurements were obtained as described below. Samples with low DNA quantity or low quality facial scans were eliminated leaving 587 DNA samples for subsequent genotyping.

Additional phenotypic trait information such as height, weight, age, sex, self-reported ancestry (based on the grandparents from both sides), eye lid (single or double), ear lobe (attached or detached), hair texture (straight, wavy, curly or very curly), freckling (none, light, medium or extensive), moles (none, few or many), as well as eye skin, and hair pigmentation was collected by a single examiner in order to reduce potential variation. The pigmentation traits were arbitrarily assigned according to previously published colour charts [26-28].

### **3D images collection and analysis**

Craniofacial scans were obtained using the Vivid 910 3-D digitiser (Konica Minolta, Australia) equipped with a medium range lens with a focal length of 14.5 mm. The scanner output images were of 640 x 480 pixels resolution for 3D and RGB data. Two daylight fluorescent sources (3400K/5400K colour temperature) were mounted at approximately 1.2 meters from the subject's head to produce ambient light conditions.

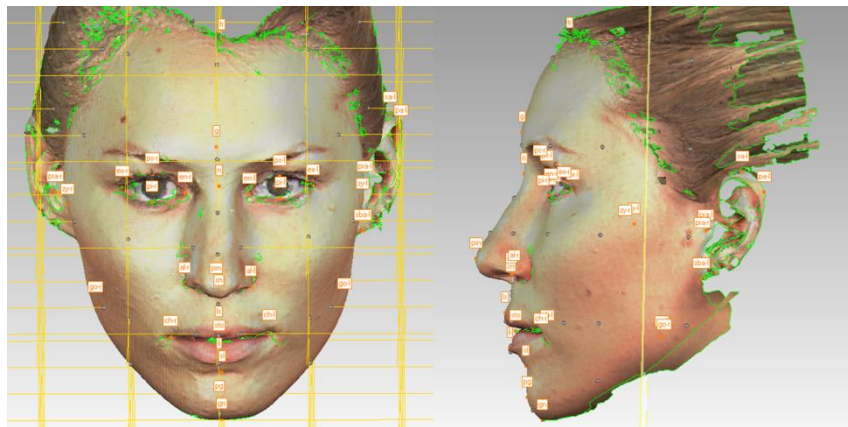
The scanner was mounted approximately one meter from the volunteer's head. Each volunteer remained in an upright seated position and kept a neutral facial expression during the scan. Subjects with long hair pulled their hair behind the ears or were asked to wear a hair net. Glasses and earrings were removed.

Each volunteer was scanned from a distance of approximately one meter from three different angles (front and two sides). The final merged 3D image was produced by semi-automatically aligning the three scans and manually cropping non-overlapping or superfluous data such as the neck area and hair using Polygone® software (Qubic, Australia). The complete coordinates of each merged 3D image were then saved in a 'vivid' file format (vvd) and exported to Geomagic® software (Qubic, Australia) for subsequent image processing.

Based on the anthropometrical literature [29] 32 anthropometrical landmarks were manually identified on each 3-D image using the Geomagic software (Fig. 1 and Supplemental Table S1). Each landmark was represented by 'x', 'y' and 'z' coordinates as part of the Cartesian coordinate system. The coordinates were exported to an Excel spreadsheet for subsequent

calculation of 86 Euclidean distances, including 54 linear distances, 10 angular distances and 21 indices (ratios) between the linear distances (Fig. 2 and Table 1).

Additionally, three direct cranial measurements: maximum cranial breadth (Euryon – Euryon), maximum cranial length (Gonion – Opisthocranium) and maximum cranial height (Vertex – Gnathion), were collected manually using a digital spreading calliper (Paleo-Tech Concepts, USA). Based on the craniofacial and body height measurements, three craniofacial ratios were calculated: Cephalic index:  $(eu-eu)/(g-op)$ , Head width – Craniofacial height index:  $(eu-eu)/(v-gn)$  and Head – Body height index:  $(v-gn)/(body\ height)$ , as summarised in Table 1.



**Figure 1. Anatomical position of the 32 manually annotated anthropometric landmarks used for calculation of linear and angular distances and ratios between the linear distances.** Some landmarks are not clearly visible due to image orientation. gn = Gnathion, pg= Pogonion, sl = Sublabiale, li = Labiale Inferius, sto = Stomion, ls = Labiale superius, ch-r = Chelion right, ch-l = Chelion left, go-r = Gonion Right, go-l = Gonion left, sn = Subnasale, prn= Pronasale, al-r = Alare right; al-l= Alare left, n = Nasion, g= Glabella; tr = Tragon, en-l = left Endocanthion, en-r = right Endocanthion, ex-r = Right Endocanthion; ex-l = left Endocanthion, ps-r = Palpebrale superius right, ps-l = Palpebrale superius left, pi-r = Palpebrale inferius right, pi-l = Palpebrale inferius left, zy-r = Zygion Right, zy-l = Zygion Left, pra-r = Tragon right, pra-l = Tragon Left, sba-l = Subalare left, sa-l = Superaurale Left, pa-l = Postaurale left.



## DNA extraction and quantification

DNA was purified from buccal swabs using the Isohelix DDK isolation kit (Cell Projects, Kent, UK) according to the manufacturer instructions. DNA samples were quantified using a Real Time quantitative PCR (q-PCR) method using a Bio-Rad CFX96 (Bio-Rad, Gladesville, Australia). This assay amplified a 63bp region of the OCA locus. The primer sequences were 5'-GCTGCAGGAGTCAGAAGGTT-3' (forward primer) and 5'-CATTTGGCGAGCAGAATCC-3' (reverse primer) at a final concentration of 200mM. All DNA samples were additionally quantified using the Qubit 2.0 fluorimeter (Invitrogen) prior to library construction as per manufacturer recommendations.

## Candidate genes and SNPs selection

Two main complementary strategies were used to generate a preliminary list of candidate genes and genetic markers. The first focused on searching the literature and web resources for candidate genes involved either in normal craniofacial variation or in craniofacial malformations in humans and model organisms (Supplemental Table S2).

The search for candidate genes focused not only on specifically defined craniofacial disorders, but also on genetic syndromes with various manifestations of craniofacial malformations, such as Down syndrome, Noonan Syndrome, Floating-Harbor Syndrome and others, as detailed in Supplemental Table S2. The main resources for locating candidate genes in the animal models were Mouse Genome Informatics [30] and AmiGo tool [31]. The main resources for identifying candidate genes in the human genome were OMIM [32] and GeneCards [33]. A comprehensive list of web resources used for candidate gene search is detailed in the Supplemental Appendix S1.

The second approach initially implemented a broad search for high Fst SNPs, such as ancestry informative markers (AIMs), with the rationale that many genes affecting craniofacial traits would have significantly different allele frequencies across populations. AIMs were selected from a variety of published and online resources [34-43].

The relevant genes obtained by both approaches were subsequently checked for potential involvement in craniofacial embryogenesis, limb development and bilateral body symmetry. It should be noted however, that the final candidate gene list was not limited to craniofacial



genes and included high  $F_{st}$  SNPs in genes with unknown function as well as markers located in intergenic regions, potentially possessing regulatory functions.

The resulting set of SNPs was further screened for high  $F_{st}$  SNPs ( $\geq 0.45$ ) in three ‘1000 genomes’ populations (CAU, ASW, CHB) using ENIGMA browser [44] as well as potentially functional polymorphisms, such as non-synonymous SNPs [45], markers in transcription factor binding sites [46] and splicing sites [47] using various web resources, as detailed in Supplemental Appendix S1 and reviewed on the GenEpi website [48]. The candidate markers search resulted in identification of 1,319 SNPs, located in approximately 177 genes/intergenic regions, as discussed in the Results section.

The chromosomal locations of final candidate markers were submitted to the custom Ampliseq primer design pipeline (Life Technologies), according to manufacturer recommendations. There were primer design difficulties for 881 markers. The marker list was therefore redesigned to include alternative tagging markers showing high linkage disequilibrium with the markers that failed initial primer design, resulting in 1,670 candidate genetic markers. Inclusion of SNPs with  $MAF < 1\%$  added an additional 4,381 genetic markers (6,051 in total). The final custom Ampliseq panel was manufactured as two separate pools of 849 and 847 primer pairs, with each amplicon covering between 125 bp and 225 bp, therefore possibly containing more than one polymorphism, and in total covering 15.78 kb of the genome. This panel included 1,319 initially targeted craniofacial and pigmentation candidate markers as well as 4,732 markers in LD with SNPs that failed primer design.

Inclusion of novel, rare SNPs ( $MAF < 1\%$ ) increased the final number of genotyped markers to 8,518 SNP in all sequenced DNA samples, although the markers with  $MAF \leq 2\%$  were not included in the association study. The list of all genotyped markers and their respective genes is detailed in Table S1.

## SNP genotyping and data analysis

Multiple DNA libraries were constructed from sets of 32 Ion Xpress<sup>TM</sup> (Life Technologies) barcoded samples using the Ion AmpliSeq<sup>TM</sup> library Kit 2.0 (Life Technologies) in conjunction with two custom primer mixes that were pooled according to manufacturer recommendations. Libraries were quantified using the Ion Library Quantitation kit (Life Technologies) and pooled in equal amounts for emulsion PCR, which was

performed using the OneTouch<sup>TM</sup> 2 instrument (Life Technologies) according to manufacturer recommendations. 587 DNA samples were genotyped by massively parallel sequencing on the Personal Genome Machine (PGM) (Life Technologies) using the Sequencing 200 v2 kit and 316 Ion chips (Life Technologies).

Raw sequencing data were collected and processed on the Torrent Suite Server v3.6.2 using default settings. Alignment and variant calling were performed against the human genome reference (hg19) sequence at low stringency settings. Binary alignment map (BAM) files were generated and exported to the Ion Reporter<sup>TM</sup> (IR) cloud-based software for SNP annotation against the reference hotspot file. The IR analysis resulted in generation of the individual variant caller files (VCF) with genotype calls for each sample as well as various statistics of the sequencing quality.

To reduce potential bias of the self-reported ancestry, ancestry inferences were obtained by 3,302 markers using STRUCTURE version 2.3.4 with default parameters as per software developer recommendations [49]. SNPs in long-range Linkage Disequilibrium (> 100,000 bp) were excluded from the STRUCTURE run. The ancestry was estimated based on four predefined population clusters: Europeans, East Asians, South Asians and Africans, according to software developer recommendations. Relative allele calls for four predefined HapMap population clusters (CEU, YRI, CHB and JPT) were used as reference populations [50]. The ancestry origin was estimated as a single (unmixed) source where the main ancestry cluster could be affiliated with at least 80% of the total mixed ancestry. The samples with mixed ancestry (>20% admixture) were assigned an ‘Admixture’ cluster.

Association analyses were performed using SNP & Variation Suite v7 (SVS) (Golden Helix, Inc., Bozeman, MT) and replicated using PLINK v1.07 software [51]. Statistical analyses in both software programs were performed using linear regression with quantitative phenotypes, and logistic regressions with binary phenotypes under the assumption of an additive genetic model, while each genotype was numerically encoded as 0, 1 or 2. Population stratification correction, incorporated by EIGENSTRAT or STRUCTURE programs was implemented in the analyses [52, 53]. In order to reduce any potential confounding effects, all the craniofacial traits association analyses were performed using sex, BMI and STRUCTURE or EIGENSTRAT ancestry clusters as covariates. In PLINK, p-values were adjusted using the ‘–adjust’ option. The final association results are based on the PLINK statistical analyses with the STRUCTURE population clusters as covariates.

Annotation analysis of the significantly associated genes was performed using the GeneCards, ENTREZ and UniProtKB web portals [33, 54]. The MalaCards web site was used to detect association between the genes and hereditary syndromes [55]. The GeneMania web site was used to identify a functional network among the genes and encoded proteins [56]. Gene ontology web resource was used to find orthologs of human genes in other organisms [31, 57]. The MGI database was used to search for the phenotype in relevant craniofacial mouse gene mutants [30]. The dbSNP, 1000 genomes, SNPnexus and Alfred websites were used for SNP annotations [58-61].

The SNP Annotation and Proxy Search (SNAP) web portal was used to find SNPs in linkage disequilibrium (LD) and generate LD plots, based on the CEU population panel from the 1000 genomes data set, within a distance of up to 500kb and an  $r^2$  threshold of 0.8 [62].

The Regulome database and potentially functional database (PFS) searches were implemented to annotate SNPs with known and predicted regulatory elements in the intergenic regions of the *H. sapiens* genome [47, 63].

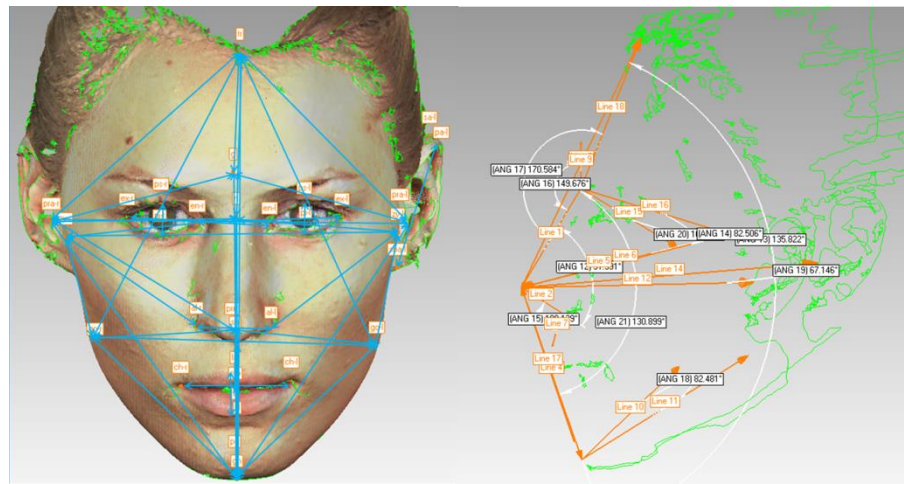
## Results and Discussion

### Phenotypic traits summary

A total of 54 linear distances, 10 angular distances and 21 indices (ratios) between the linear distances were calculated based on the Cartesian coordinates of 32 anthropometric landmarks that were manually mapped on each of the 587 3-D facial images (Fig. 1, Fig. 2, Table 1 and Supplemental Table S1). Three additional craniofacial distances were obtained by direct measurement of subjects' heads and used to calculate three indices: maximum cranial breadth, maximum cranial length and maximum cranial height, cephalic index, head width – craniofacial height index and head – body height index (Table 1). Furthermore, the linear and angular facial distances were used to calculate 20 principal components (PCs). Additional phenotypic features such as eyelid, pigmentation, hair texture, freckling, moles, height, weight, BMI, age, sex and ancestry were collected. In total, the data on 104

phenotypic traits were recorded and used for genetic association analyses, although this study focuses only on the anthropometric craniofacial phenotypes association.

The phenotypic data collection by a single examiner achieved more consistent measurements from the 3-D image analyses. In addition, all measurements were based on the images of participants within a narrow age range 26.6 (SD  $\pm$  8.9)



**Figure 2. Illustration of linear and angular distances calculated from manually annotated landmark coordinates.**

**Table 1. Craniofacial anthropometric measurements that were calculated and used for genetic association analyses.**

**Manual craniofacial measurements**

- V-Gn (Maximum Craniofacial height)
- Eu-Eu (Maximum Head Width)
- G-Op (Maximum Head Length)
- Cephalic index: (eu-eu)/(g-op)
- Head width – Craniofacial height index: (eu-eu)/(v-gn)
- Head – Body height index: (v-gn)/(body height)

**3D facial measurements**

**Linear facial distances**

- Total face height: tr-gn
- Face width: zy-zy

- Morphological face height: n-gn
- Physiognomical face height: n-sto
- Lower profile height: prn-gn
- Lower face height: sn-gn
- Lower third face depth: t(l)-gn
- Middle face depth: t(l)-prn
- Middle face height (right): go(r)-zy(r)
- Middle face height (left): go(l)-zy(l)
- Middle face width 1: t(r)-t(l)
- Middle face width 2 (left): zy(l)-al(l)
- Middle face width 2 (right): zy(r)-al(r)
- Upper face depth: (left): t(l)-tr
- Upper face depth: (right): t(r)-tr
- Upper third face depth: t(l)-n
- Forehead height: g-tr
- Extended forehead height: tr-n
- Glabella –Gnathion distance: g-gn
- Supraorbital depth: t(l)-g
- Trichion – Zygion distance (left): tr-zy(l)
- Trichion – Zygion distance (right): tr-zy(r)
- Nasion - Zygion distance (left): n-zy(l)
- Nasion - Zygion distance (right): n-zy(r)
- Zygion – Gnathion distance (left): zy(l)-gn
- Zygion – Gnathion distance (right): zy(r)-gn
- Interendocanthal width: en-en
- Interexocanthal width: ex-ex
- Eye fissure width (left): en(l)-ex(l)
- Eye fissure width (right): en(r)-ex(r)
- Eye fissure height (left): ps(l)-pi(l)
- Eye fissure height (right): ps(r)-pi(r)
- Ear height (left): sa(l)-sba(l)
- Ear width (left): t(l)-pa(l)
- Nasal bridge length: n-prn
- Nose height: n-sn
- Nose width: al-al
- Nasal tip protrusion: sn-prn
- Ala length (left): prn-al(l)
- Ala length (right): prn-al(r)
- Gonion - Trichion distance (left): go(l)-tr
- Gonion - Trichion distance (right): go(r)-tr
- Gonion – Glabella distance: g-pg

- Pronasale - Gonion distance (left): prn-go(l)
- Pronasale - Gonion distance (right): prn-go(r)
- Chin height: sl-gn
- Mandibular region depth (right): t(r)-gn
- Mandible width: go-go
- Mandible height: sto-gn
- Lower jaw depth (left): gn-go(l)
- Lower jaw depth (right): gn-go(r)
- Mouth width: ch-ch
- Upper vermilion height: ls-sto
- Lower vermilion height: li-sto

### **Angular facial distances**

- Nasal tip angle: (n-prn-sn)
- Nasal vertical prominence angle: (tr-prn-gn)
- Transverse nasal prominence angle 1: (zy(l)-prn-zy(r))
- Transverse nasal prominence angle 2: (t(l)-prn-t(r))
- Nasolabial angle: (prn-sn-ls)
- Nasofrontal angle: (g-n-prn)
- Nasion depth angle: (zy(l)-n-zy(r))
- Nasomental angle: (n-prn-pg)
- Forehead nasal angle: (tr-n-prn)
- Chin prominence angle: (go(l)-gn-go(r))

### **Ratios (indices)**

- Forehead height ratio:  $(tr-n)/(go(r)-go(l))$
- Upper face height ratio:  $(n-sn)/(go(r)-go(l))$
- Lower face height ratio:  $(sn-gnx)/(go-go)$
- Anterior face height 1 ratio:  $(n-gn)/(go-go)$
- Anterior face height 2 ratio:  $(n-gn)/(zy-zy)$
- Face height index:  $(n-gn)/(tr-gn)$
- Upper – Lower face ratio:  $(tr-g)/(sn-gn)$
- Upper face height ratio:  $(n-sn)/(sn-gn)$
- Upper face width ratio:  $(n-sn)/(zy-zy)$
- Total anterior face height ratio:  $(tr-gn)/(zy-zy)$
- Mouth width ratio:  $(ch-ch) \times 100 / (en-en)$
- Mandible – Face width ratio:  $(go-go)/(zy-zy)$
- Mandible index:  $(sto-gn) \times 100 / (go-go)$
- Mandible – Interexocanthion distance ratio:  $(go-go)/(ex-ex)$
- Interendocanthion distance ratio:  $(en-en)/(al-al)$

- Intercanthal index:  $(en(r)-en(l))/(ex(r)-ex(l))$
- Intercanthal – Intracanthal index:  $(ex(r)-en(r))/(en(l)-ex(l))$
- Nasal index:  $(al-al) \times 100 / (n-sn)$
- Nose-face height index:  $(n-sn) / (n-gn)$
- Nose-face width index:  $(al-al) / (zy-zy)$
- Nasal tip protrusion – nose width index:  $(sn-prn) / (al-al)$
- Nasal tip protrusion –Nose height index:  $(sn-prn) / (n-sn)$

### 3D measurements precision study

In the last decade 3D scanning systems have been extensively used in anthropometric studies as well as in medical research [18, 20, 64, 65]. The Minolta Vivid V910 3D scanner has been demonstrated to have accuracy to a level of  $1.9 \pm 0.8$  mm [66] and  $0.56 \pm 0.25$  mm [67], making it suitable for the present study since it should provide an accurate representation of facial morphology. However, the allocation of anthropometric facial landmarks can be challenging, especially when tissue palpating is not possible. Reproducibility of the landmark precision was assessed on fifteen 3D facial images through assessment of 86 facial measurements, including linear and angular distances and ratios between the linear distances at two separate times. The period between the analyses varied from one to six months. The mean difference (MD) was calculated as the discrepancy between the first and the second measurement. The measurement error (ME) was calculated as the standard deviation of the MD divided by square root of 2 ( $ME=SD(MD/\sqrt{2})$ ).

In general, the nasal area distances, which involved nasion, pronasale, subnasale and alare landmarks showed greater reproducibility, while the measurements involving paired landmarks, such as gonion and zygion demonstrated higher variance. This result can be explained by easier allocation of nasal area landmarks, compared with gonion and zygion [29]. Overall the median difference (MD) between two measurements for linear distances in 15 images ranged between 0.76 mm ( $ME \pm 0.27$ ) and 2.80 mm ( $ME \pm 0.99$ ); for angular distances between 0.38 mm ( $ME \pm 0.96$ ) and 3.75 mm ( $ME \pm 0.40$ ) and for facial indices (ratios) between 0.46 mm ( $ME \pm 1.08$ ) and 2.98 mm ( $ME \pm 1.95$ ) respectively. The lower reproducibility in the angular distances and indices can be explained by a higher number of landmarks (hence variability in allocation of x, y and z coordinates) needed for their calculation (three and four landmarks respectively). Nevertheless, our findings are concordant



with the published results, which observed variance of 0.19 mm to 3.49 mm with a ME range of 0.55 mm to 3.34 mm for each landmark [19, 68].

## **Candidate genes search and sequencing data quality control**

The search for candidate genes and SNPs potentially involved in influencing normal craniofacial morphology variation initially focused on searching for genes involved in normal or abnormal craniofacial variation in humans and model organisms (Supplemental Table S2). As a complementary approach, a search for genetic markers with high  $F_{st}$  values ( $\geq 0.45$ ) was implemented, based on the rationale that genes involved in craniofacial morphology regulation are likely to display significant differences in allele frequencies across populations.

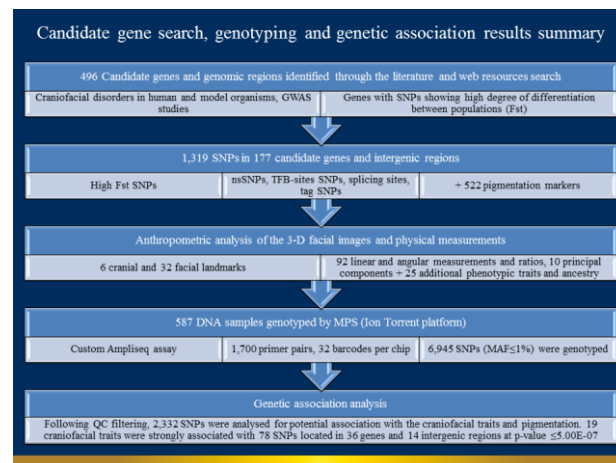
The first approach has mainly focused on the Mouse Genome Informatics (MGI) database search using the keyword ‘craniofacial mutants’ and additional resources such as Online Mendelian Inheritance in Man (OMIM), GeneCards and AmiGO, using the keywords such as “craniofacial”, “craniofacial mutants”, “craniofacial anomalies”, “craniofacial dimorphism” and “facial morphology” (a detailed list of used resources is summarized in Supplemental Appendix S1). This search revealed a list of 2,891 genotypes and 7,956 annotations. A search of the ‘abnormal facial morphology’ sub-category resulted in 1,492 genotypes and 2,889 annotations. The final search of the ‘abnormal nose morphology’ of the previous sub-category revealed 219 genotypes with 310 annotations, representing approximately 150 genes.

In parallel, a search high  $F_{st}$  markers, using previously published AIMs and web tools, such as ENGINES, resulted in identification of additional targets, for a total of 1,088 genes and intergenic regions (a detailed list of used resources is summarized in Supplemental Appendix S1).

However, manual examination revealed that 592 of these genes showed no apparent link with normal craniofacial development or malformations and were therefore excluded. The remaining 496 regions were further screened for non-synonymous and potentially functional SNPs, as well as SNPs with high population differentiation, which resulted in the shortlist of 269 genes and intergenic regions.

Subsequent analysis of these 269 genes/regions for functional annotation using the AmiGO Gene Ontology server [57], resulted in 177 candidate genes/regions, possessing 1,319

genetic markers involved in various stages of human embryonic development, including: embryonic morphogenesis, sensory organ development, tissue development, pattern specification process, tissue morphogenesis, ear development, tube morphogenesis, epithelium development, chordate embryonic development and morphogenesis of an epithelium (Supplemental Appendix S1). Notably, the majority of these markers are located in introns and intergenic regions (summarized in Fig. 3).



**Figure 3. A simplified flow chart, summarizing major steps of the candidate gene search, SNP genotyping and data analysis workflow.**

In terms of molecular function, AmiGO showed that craniofacial candidate markers might be involved in a range of regulatory activities including: protein dimerization activity, chromatin binding, regulatory region DNA binding, sequence-specific DNA binding RNA polymerase II transcription factor activity, sequence-specific distal enhancer binding activity, heparin binding, RNA polymerase II core promoter proximal region sequence-specific DNA binding transcription factor activity involved in positive regulation of transcription, BMP receptor binding and transmembrane receptor protein serine/threonine kinase binding (Supplemental Appendix S1).

Subsequent analysis of candidate SNPs for mouse phenotype associations confirmed that orthologous candidate markers were previously detected in mouse models displaying abnormal morphology of the skeleton, head, viscerocranium and facial area, as well as specific malformations of the eye, ear, jaw, palate, limbs, digits and tail (data not shown).

In addition to craniofacial candidate SNPs, 522 markers, previously shown to be associated with pigmentation traits, such as eye, skin and hair colour were selected from the relevant literature. These markers were used to validate the results of the genetic association analyses of craniofacial traits.

The final candidate marker list was analysed using the GREAT platform to visualize the genomic context of amplicons covering targeted SNPs [69]. The analysis revealed that almost 99% of the genomic regions (which may cover multiple markers) are associated with one or two genes with approximately 62% of genomic regions located 0-500 kb downstream of a transcription start site (data not shown).

Targeted mass parallel sequencing of the 587 samples resulted in 9,051 genetic markers, with the majority of markers (>5,000) represented by rare polymorphisms of  $\leq 1\%$  minor allele frequency (MAF) (data not shown). The difference between the initial hot-spot SNP panel of candidate markers ( $n=6,945$ ) and the actual sequencing output ( $n=9,051$ ) was a result of identification of potentially novel and rare markers in individual DNA samples. Three of the 587 samples, did not produce high quality genotypes because of poor DNA quality or unsuccessful library and template preparation (summarized in Fig. 3).

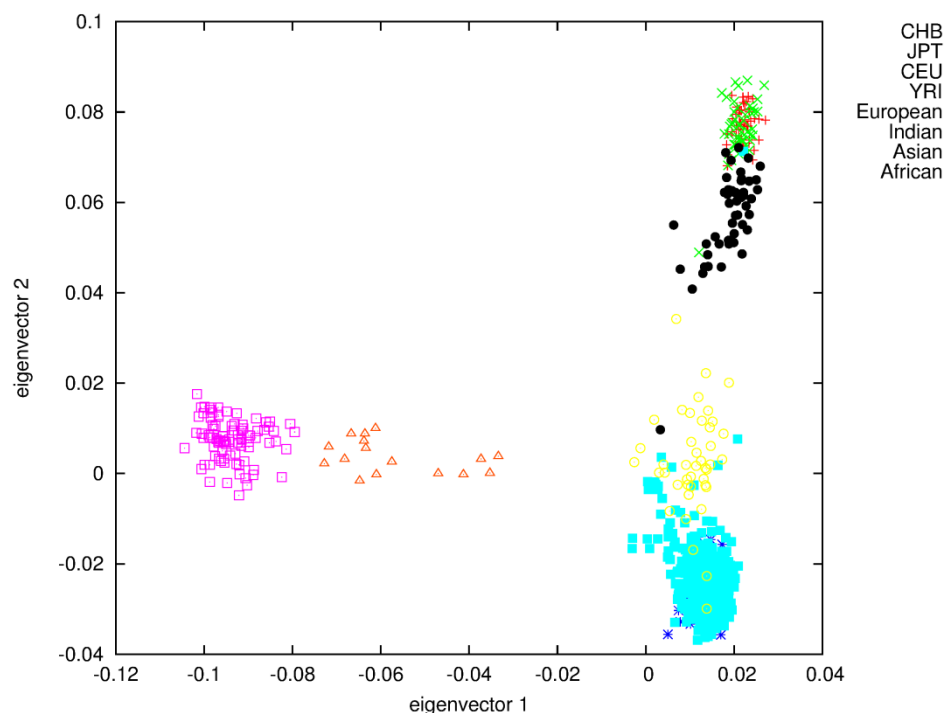
The SNPs were filtered by sequencing quality and by MAF. Data quality control was performed by removing markers of low genotype quality ( $GQ > 10$ ) and sequencing depth ( $DP > 10X$ ), which resulted in 8,518 markers (Supplemental Appendix S2). Further filtering of markers using a 2% MAF cut-off resulted in 3,073 markers (Supplemental Appendix S2). The decision to apply a slightly more stringent MAF threshold (2%) was made because of the sample size ( $n=587$ ) and to reduce potential bias from rare SNPs (1% MAF). Since this may reduce the power of analysis, we analysed and compared both datasets and did not observed any significant difference. Additional filtering based on the HWE threshold of  $p\text{-value} \geq 0.01$  resulted in 3,073 markers. The mean sequencing depth for significantly associated markers in this study was 58 fold ( $\pm 48.9$  SD).

## Genetic association study

The association analyses were performed using a linear regression model, incorporating a covariate for potential population stratification as well as sex and BMI as covariates. The use of covariates in the statistical analysis aimed to reduce the risk of

introducing confounding effects, which can result in false positive associations. While sexual dimorphism in the craniofacial morphology is well-known [70], BMI will also likely affect certain craniofacial traits, since the soft facial tissue may change significantly with weight gain or loss. Despite that, this potential confounding factor has to date been disregarded in association studies of normal craniofacial morphology. Age was not considered a significant covariate, given that average age of the subjects in this study was 27 ( $\pm 8.9$  SD). Nevertheless, the potential effect of age as a cofactor was assessed on three craniofacial traits and found to be not significant (data not shown).

In contrast to most other craniofacial association studies that focused on a specific homogeneous population group (mostly Europeans), this study included samples from several population groups, which enabled investigation of the genetic factors influencing normal craniofacial morphology in different ethnicities [71]. Self-reported ancestry however, cannot be considered fully reliable, as demonstrated previously [72, 73]. In order to address this issue we assessed the self-reported ancestry using STRUCTURE with 186 SNPs removed due to long-range disequilibrium [49]. Following the rationale that the best ancestry estimates are obtained using a large number of random markers [74], we used all the available markers (after MAF filtering) in STRUCTURE analysis. The STRUCTURE analysis resulted in clusters of 367 Europeans, 51 East Asians, 43 South Asians and 16 Africans, with 107 samples designated as admixed ancestry (Fig. 4). Of the samples tested with STRUCTURE, 459 (89%) were assigned the same ancestry cluster (sole or mixed origin) as the self-reported information. Of the remaining 57 individuals, 39 were estimated as ‘admixture’ (based on up to 20% admixture threshold) and 18 were assigned a single ancestry, different to the self-reported ancestry (Fig. 4).



**Figure 4. Population structure as represented by plotting genomic PCs 1 and 2, using 270 HapMap individuals.** YRI: Yoruba, Nigeria, Africa. JRI: Japanese, Tokyo, Japan. CHB: Han Chinese, Beijing, China. CEU: Utah residents with European ancestry.

The risk of detecting false positive results because of population stratification was carefully assessed and further reduced by applying an EIGENSTRAT correction. Specifically, EIGENSTRAT's smartpca.perl was used to perform PCA-clustering in comparison to reference populations from HapMap reference clusters. The Q-Q plots of the associated traits showed the expected distribution of data after the applied correction (Supplemental Figs S1-S14).

We did not perform allele imputations on this dataset because it includes individuals from heterogenous ancestral backgrounds, with 107 subjects classified as 'admixture', based on the applied threshold of 20%. Imputation using homogenous reference populations would have introduced unnecessary bias with wrongly imputed alleles in subsequent analysis steps.

While the majority of current GWA studies rely on a p-value  $<5.00E-08$  significance threshold, some publications suggest this threshold may be too stringent, especially for

complex traits that are regulated by a large number of small effect alleles [75, 76]. In contrast to GWAS, candidate gene studies undertake a more focused genetic strategy, concentrating on a relatively limited number of putative markers. As this study analysed a much smaller number of SNPs than usual GWA-studies, we could use a higher p-value cut-off since the smaller sample size means the probability of false positive at extremely low p-values is itself lower. Nevertheless, we decided to keep the traditional GWAS p-value significance threshold ( $<5.00\text{E-}08$ ) in order to reduce the possibility of detecting false positive results.

In addition, we subsequently applied a more stringent Bonferroni –corrected threshold in order to minimize the chance of detecting spurious associations. Following the association analysis of 104 craniofacial phenotypes with 3,073 genetic markers, the significance threshold based on the Bonferroni correction with a desired  $\alpha$  of 0.05 would be  $1.6\text{E-}07$  ( $=0.05/(3073*104)$ ).

However, it should be emphasized that the Bonferroni correction is widely considered over-conservative, especially in the case of complex phenotypic traits with small individual effects of each allele. Considering that our results confirm the previously published findings, we believe the GWAS p-value threshold is conservative enough to avoid or at least significantly reduce potentially spurious associations. Following this rationale, we report all the variants, which met the GWAS (and subsequently the Bonferroni) p-value threshold.

In our attempt to identify genetic markers influencing normal variation in craniofacial traits, we incorporated 522 markers previously associated with human pigmentation traits, such as eye, skin and hair colour. These markers were included to validate the statistical methods used for the craniofacial traits association study. The association analyses of the pigmentation traits, which were based on the HWE non-filtered data, did indeed confirm previously published findings, as detailed in Table S3. It should be noted however, that these results may not necessarily confirm the validity of the craniofacial markers associations.

The application of the Hardy-Weinberg equilibrium (HWE) threshold resulted in filtering 25% of the total number of SNPs. These markers included almost all the SNPs, previously associated with pigmentation traits, such as rs12913832, rs1129038, rs8039195 and rs16891982. This is not surprising, since population-related markers (both pigmentation and craniofacial SNPs) are likely not being in HWE ‘a priori’. Another explanation for this observation is potential bias from partially uncorrected heterogeneous ancestry, since the ancestry correction algorithm can only minimize, rather than completely remove spurious

associations [52]. In fact, the association analyses of the HWE non-filtered genotyping data with pigmentation traits (eye, skin and hair colour), demonstrated highly significant associations, concordant with the literature (Table S3).

The results of the association analyses of the craniofacial traits are summarized in Table 2 and Supplemental Figs. S15-S28. In general, two linear distances, four angular distances, four indices and one PC revealed significant associations with 45 SNPs in 27 genes and 13 intergenic regions, based on the  $5.00\text{E-}08$  p-value threshold, including three SNPs, associated with three linear distances that showed a borderline p-values (Table 2). However, when applying a more stringent Bonferroni correction (p-value  $1.6\text{E-}07$ ), this list could be reduced to 12 SNPs in 12 genes and intergenic regions that were associated with 8 phenotypic traits (Table 2). The genes and the significant SNPs are: AGXT2 (rs37369); FOXN3 (rs390345); TEX41 (rs10496971); PCDH15 (rs10825273); DCT (rs1407995); CACNB4 (rs16830498); ZEB1 (rs59037879); FAM49A (rs6741412); STON1 (rs7574549); and three intergenic SNPs (rs10512572, rs8035124, rs942316). All the SNP are located in introns, except rs37369, which represents a missense mutation.



**Table 2. Results of genetic association analyses between candidate SNPs and craniofacial traits**, including all the genomic markers that reached the GWAS p-value threshold ( $5.00E-08$ ). Highlighted with **bold**: genomic markers that reached Bonferroni corrected threshold ( $1.6E-07$ ); Highlighted with blue colour: linear distances; Highlighted with red colour: craniofacial indices; Highlighted with green colour: angular distances; Highlighted with violet colour: principal component; Gene: gene name; rs#: reference SNP ID number; SNP: chromosomal location of the marker; Genomic annotation: genomic location of the marker; UNADJ: Unadjusted p-values. BONF: Bonferroni single-step adjusted. HOLM: Holm (1979) step-down adjusted. SIDAK\_SS: Sidak single-step adjusted. SIDAK\_SD: Sidak step-down adjusted. FDR\_BH: Benjamini & Hochberg (1995) step-up FDR control. FDR\_BY: Benjamini & Yekutieli (2001) step-up FDR control.

Gene	rs#	SNP	Genomic annotation	UNADJ	BONF	HOLM	SIDAK_SS	SIDAK_SD	FDR_BH	FDR_BY
<b>al-al</b>										
<b>between SV2B and TRNAY16P</b>	<b>rs8035124</b>	<b>15:92105708</b>	<b>intergenic</b>	<b>1.52E-10</b>	<b>1.74E-07</b>	<b>1.74E-07</b>	<b>1.74E-07</b>	<b>1.74E-07</b>	<b>1.74E-07</b>	<b>1.33E-06</b>
EYA2	rs58733120	20:45803852	intronic	5.37E-10	6.15E-07	6.14E-07	6.15E-07	6.14E-07	2.93E-07	2.23E-06
RP11-494M8.4	rs1482795	11:7850345	intergenic	7.68E-10	8.78E-07	8.77E-07	8.78E-07	8.77E-07	2.93E-07	2.23E-06
AGXT2	rs37369	5:35037115	missense	1.04E-09	1.19E-06	1.19E-06	1.19E-06	1.19E-06	2.98E-07	2.27E-06
downstream to PTCH1	rs57585041	9:98205221	intergenic	6.05E-09	6.92E-06	6.90E-06	6.92E-06	6.90E-06	1.38E-06	1.06E-05
EYA1	rs79867447	8:72127562	intronic	3.92E-08	4.49E-05	4.47E-05	4.49E-05	4.47E-05	7.48E-06	5.70E-05
<b>sn-prn</b>										
between LOC100131241 and LOC124685	rs10512572	17:69512099	intergenic	2.22E-08	2.54E-05	2.54E-05	2.54E-05	2.54E-05	2.54E-05	1.94E-04
<b>cephalic index</b>										
<b>CACNB4</b>	<b>rs16830498</b>	<b>2:152814028</b>	<b>intronic</b>	<b>7.57E-11</b>	<b>8.67E-08</b>	<b>8.67E-08</b>	<b>8.67E-08</b>	<b>8.67E-08</b>	<b>8.67E-08</b>	<b>6.61E-07</b>
MYO5A	rs2290332	15:52611451	synonymous	5.56E-10	6.37E-07	6.37E-07	6.37E-07	6.37E-07	2.39E-07	1.82E-06
ZEB1	rs59037879	10:31745993	intronic	6.27E-10	7.18E-07	7.17E-07	7.18E-07	7.17E-07	2.39E-07	1.82E-06
COL11A1	rs4908280	1:103420759	intronic	1.66E-09	1.91E-06	1.90E-06	1.91E-06	1.90E-06	4.73E-07	3.61E-06

EYA1	rs1481800	8:72131426	intronic	2.07E-09	2.37E-06	2.36E-06	2.37E-06	2.36E-06	4.73E-07	3.61E-06
TEX41	rs10496971	2:145769943	intronic	5.32E-09	6.09E-06	6.06E-06	6.09E-06	6.06E-06	1.01E-06	7.73E-06
PCDH15	rs10825273	10:55968685	intronic	9.93E-09	1.14E-05	1.13E-05	1.14E-05	1.13E-05	1.62E-06	1.24E-05
COL11A1	rs11164649	1:103444679	intronic	1.70E-08	1.95E-05	1.94E-05	1.95E-05	1.94E-05	2.43E-06	1.86E-05
-	rs373272	5:84818656	intergenic	2.40E-08	2.75E-05	2.73E-05	2.75E-05	2.73E-05	3.06E-06	2.33E-05
<b>nasal index</b>										
<b>AGXT2</b>	<b>rs37369</b>	<b>5:35037115</b>	<b>missense</b>	<b>1.25E-10</b>	<b>1.43E-07</b>	<b>1.43E-07</b>	<b>1.43E-07</b>	<b>1.43E-07</b>	<b>1.43E-07</b>	<b>1.09E-06</b>
EYA2	rs58733120	20:45803852	intronic	9.46E-09	1.08E-05	1.08E-05	1.08E-05	1.08E-05	5.23E-06	3.99E-05
RP11-408B11.2	rs7311798	12:85808703	intergenic	1.77E-08	2.02E-05	2.02E-05	2.02E-05	2.02E-05	5.23E-06	3.99E-05
-	rs1482795	11:7850345	intergenic	1.83E-08	2.09E-05	2.09E-05	2.09E-05	2.09E-05	5.23E-06	3.99E-05
EYA1	rs79867447	8:72127562	intronic	3.53E-08	4.03E-05	4.02E-05	4.03E-05	4.02E-05	8.07E-06	6.15E-05
<b>nasal tip protrusion-width index</b>										
<b>AGXT2</b>	<b>rs37369</b>	<b>5:35037115</b>	<b>missense</b>	<b>1.09E-11</b>	<b>1.24E-08</b>	<b>1.24E-08</b>	<b>1.24E-08</b>	<b>1.24E-08</b>	<b>1.24E-08</b>	<b>9.46E-08</b>
<b>between LOC100131241 and LOC124685</b>	<b>rs10512572</b>	<b>17:69512099</b>	<b>intergenic</b>	<b>4.48E-11</b>	<b>5.12E-08</b>	<b>5.12E-08</b>	<b>5.12E-08</b>	<b>5.12E-08</b>	<b>2.56E-08</b>	<b>1.95E-07</b>
<b>TEX41</b>	<b>rs10496971</b>	<b>2:145769943</b>	<b>intronic</b>	<b>9.36E-11</b>	<b>1.07E-07</b>	<b>1.07E-07</b>	<b>1.07E-07</b>	<b>1.07E-07</b>	<b>3.57E-08</b>	<b>2.72E-07</b>
FOXN3	rs390345	14:89976534	intronic	1.96E-10	2.24E-07	2.23E-07	2.24E-07	2.23E-07	5.60E-08	4.27E-07
EYA2	rs58733120	20:45803852	intronic	9.28E-10	1.06E-06	1.06E-06	1.06E-06	1.06E-06	2.12E-07	1.62E-06
LMNA	rs12076700	1:156055099	intronic	1.66E-09	1.90E-06	1.89E-06	1.90E-06	1.89E-06	2.81E-07	2.14E-06
FAM49A	rs11096686	2:16815892	intronic	1.72E-09	1.97E-06	1.96E-06	1.97E-06	1.96E-06	2.81E-07	2.14E-06
MEIS2	rs2122497	15:37247246	intronic	2.33E-09	2.67E-06	2.65E-06	2.67E-06	2.65E-06	3.33E-07	2.54E-06
PCDH15	rs10825273	10:55968685	intronic	2.82E-09	3.22E-06	3.20E-06	3.22E-06	3.20E-06	3.58E-07	2.73E-06
-	rs11004765	10:56918449	intergenic	1.66E-08	1.89E-05	1.88E-05	1.89E-05	1.88E-05	1.89E-06	1.44E-05
<b>nose-face width index</b>										
EYA1	rs79867447	8:72127562	intronic	1.10E-09	1.26E-06	1.26E-06	1.26E-06	1.26E-06	1.26E-06	9.62E-06
EYA2	rs58733120	20:45803852	intronic	3.38E-09	3.86E-06	3.86E-06	3.86E-06	3.86E-06	1.93E-06	1.47E-05

EGFR	rs17335905	7:55131384	intronic	4.74E-08	5.42E-05	5.41E-05	5.42E-05	5.41E-05	1.81E-05	1.38E-04
nasion depth angle										
TEX41	rs10496971	2:145769943	intronic	1.86E-17	2.13E-14	2.13E-14	INF	INF	2.13E-14	1.62E-13
upstream to BMP4	rs942316	14:54440983	intergenic	3.74E-14	4.28E-11	4.28E-11	4.28E-11	4.28E-11	2.14E-11	1.63E-10
CACNB4	rs16830498	2:152814028	intronic	5.91E-14	6.76E-11	6.75E-11	6.76E-11	6.75E-11	2.25E-11	1.72E-10
PCDH15	rs10825273	10:55968685	intronic	1.55E-12	1.78E-09	1.77E-09	1.78E-09	1.77E-09	4.44E-10	3.39E-09
between LOC100131241 and LOC124685	rs10512572	17:69512099	intergenic	2.89E-12	3.30E-09	3.29E-09	3.30E-09	3.29E-09	6.61E-10	5.04E-09
FAM49A	rs6741412	2:16815759	intronic	3.00E-11	3.43E-08	3.42E-08	3.43E-08	3.42E-08	5.72E-09	4.36E-08
ZEB1	rs59037879	10:31745993	intronic	5.77E-11	6.60E-08	6.56E-08	6.60E-08	6.56E-08	9.42E-09	7.18E-08
DCT	rs1407995	13:95096013	intronic	9.42E-11	1.08E-07	1.07E-07	1.08E-07	1.07E-07	1.35E-08	1.03E-07
AGXT2	rs37369	5:35037115	missense	1.31E-10	1.50E-07	1.49E-07	1.50E-07	1.49E-07	1.67E-08	1.27E-07
ASTN2	rs10513300	9:120130206	intronic	2.37E-10	2.71E-07	2.69E-07	2.71E-07	2.69E-07	2.48E-08	1.89E-07
MBD5	rs7569399	2:149234861	intronic	2.39E-10	2.73E-07	2.71E-07	2.73E-07	2.71E-07	2.48E-08	1.89E-07
between FRMD1 and CTAGE13P	rs1871428	6:168665760	intergenic	2.70E-10	3.09E-07	3.06E-07	3.09E-07	3.06E-07	2.57E-08	1.96E-07
MYO5A	rs2290332	15:52611451	synonymous	3.67E-10	4.20E-07	4.15E-07	4.20E-07	4.15E-07	2.93E-08	2.24E-07
BCOR	rs5963728	X:39919182	intronic	3.67E-10	4.20E-07	4.15E-07	4.20E-07	4.15E-07	2.93E-08	2.24E-07
DCT	rs1325611	13:95094385	intronic	3.85E-10	4.40E-07	4.35E-07	4.40E-07	4.35E-07	2.93E-08	2.24E-07
LMNA	rs12076700	1:156055099	intronic	1.09E-09	1.24E-06	1.23E-06	1.24E-06	1.23E-06	7.77E-08	5.92E-07
STON1	rs7574549	2:48822540	intronic	1.33E-09	1.52E-06	1.50E-06	1.52E-06	1.50E-06	8.97E-08	6.83E-07
between HAS2-AS1 and MRPS36P3	rs7844723	8:122908503	intergenic	1.56E-09	1.79E-06	1.76E-06	1.79E-06	1.76E-06	9.93E-08	7.57E-07
KIAA1217	rs3910620	10:24457977	intronic	2.16E-09	2.47E-06	2.43E-06	2.47E-06	2.43E-06	1.30E-07	9.89E-07
nasolabial angle										
SMAD1	rs17020235	4:146418167	intronic	2.07E-09	2.36E-06	2.36E-06	2.36E-06	2.36E-06	2.36E-06	1.80E-05
transverse nasal prominence angle 1										
TEX41	rs10496971	2:145769943	intronic	6.80E-14	7.78E-11	7.78E-11	7.77E-11	7.77E-11	7.78E-11	5.93E-10

FAM49A	rs6741412	2:16815759	intronic	3.22E-12	3.69E-09	3.68E-09	3.69E-09	3.68E-09	1.84E-09	1.40E-08
between LOC100131241 and LOC124685	rs10512572	17:69512099	intergenic	1.50E-11	1.71E-08	1.71E-08	1.71E-08	1.71E-08	5.70E-09	4.35E-08
upstream to BMP4	rs942316	14:54440983	intergenic	5.51E-11	6.30E-08	6.28E-08	6.30E-08	6.28E-08	1.58E-08	1.20E-07
STON1	rs7574549	2:48822540	intronic	1.27E-10	1.45E-07	1.45E-07	1.45E-07	1.45E-07	2.90E-08	2.21E-07
LHX8	rs12041465	1:75609049	intronic	2.46E-10	2.82E-07	2.81E-07	2.82E-07	2.81E-07	4.70E-08	3.58E-07
PCDH15	rs10825273	10:55968685	intronic	3.36E-10	3.85E-07	3.83E-07	3.85E-07	3.83E-07	5.49E-08	4.19E-07
CACNB4	rs16830498	2:152814028	intronic	6.16E-10	7.05E-07	7.01E-07	7.05E-07	7.01E-07	8.25E-08	6.29E-07
DCT	rs1407995	13:95096013	intronic	6.49E-10	7.43E-07	7.37E-07	7.43E-07	7.37E-07	8.25E-08	6.29E-07
AGXT2	rs37369	5:35037115	missense	8.65E-10	9.90E-07	9.82E-07	9.90E-07	9.82E-07	9.90E-08	7.54E-07
LMNA	rs12076700	1:156055099	intronic	1.30E-09	1.48E-06	1.47E-06	1.48E-06	1.47E-06	1.34E-07	1.02E-06
RTTN	rs74884233	18:67813813	intronic	1.40E-09	1.60E-06	1.59E-06	1.60E-06	1.59E-06	1.34E-07	1.02E-06
DCT	rs1325611	13:95094385	intronic	2.12E-09	2.43E-06	2.40E-06	2.43E-06	2.40E-06	1.87E-07	1.42E-06
ZEB1	rs59037879	10:31745993	intronic	5.65E-09	6.47E-06	6.39E-06	6.47E-06	6.39E-06	4.62E-07	3.52E-06
KIAA1217	rs3910620	10:24457977	intronic	1.06E-08	1.21E-05	1.20E-05	1.21E-05	1.20E-05	8.09E-07	6.16E-06
BCOR	rs5963728	X:39919182	intronic	1.28E-08	1.47E-05	1.45E-05	1.47E-05	1.45E-05	9.18E-07	7.00E-06
between FRMD1 and CTAGE13P	rs1871428	6:168665760	intergenic	1.37E-08	1.56E-05	1.54E-05	1.56E-05	1.54E-05	9.20E-07	7.01E-06
TMTC2	rs11115552	12:83423379	intronic	1.97E-08	2.25E-05	2.22E-05	2.25E-05	2.22E-05	1.25E-06	9.53E-06
FGF14	rs2476230	13:102582443	intronic	3.10E-08	3.55E-05	3.50E-05	3.55E-05	3.50E-05	1.87E-06	1.42E-05
transverse nasal prominence angle 2										
ZEB1	rs59037879	10:31745993	intronic	5.31E-12	6.07E-09	6.07E-09	6.07E-09	6.07E-09	6.07E-09	4.62E-08
between LOC100131241 and LOC124685	rs10512572	17:69512099	intergenic	1.38E-11	1.57E-08	1.57E-08	1.57E-08	1.57E-08	7.85E-09	5.98E-08
AGXT2	rs37369	5:35037115	missense	1.46E-09	1.66E-06	1.66E-06	1.66E-06	1.66E-06	4.38E-07	3.34E-06
LMNA	rs12076700	1:156055099	intronic	1.54E-09	1.75E-06	1.75E-06	1.75E-06	1.75E-06	4.38E-07	3.34E-06
FAM49A	rs6741412	2:16815759	intronic	2.75E-09	3.14E-06	3.13E-06	3.14E-06	3.13E-06	6.28E-07	4.78E-06
TEX41	rs10496971	2:145769943	intronic	5.52E-09	6.30E-06	6.27E-06	6.30E-06	6.27E-06	1.05E-06	8.00E-06

RTTN	rs74884233	18:67813813	intronic	1.20E-08	1.37E-05	1.37E-05	1.37E-05	1.37E-05	1.96E-06	1.50E-05
AC073218.1	rs892458	2:34667749	intergenic	1.73E-08	1.98E-05	1.97E-05	1.98E-05	1.97E-05	2.48E-06	1.89E-05
PAX3	rs2289266	2:223089431	intronic	1.95E-08	2.23E-05	2.21E-05	2.23E-05	2.21E-05	2.48E-06	1.89E-05
LHX8	rs12041465	1:75609049	intronic	2.30E-08	2.62E-05	2.60E-05	2.62E-05	2.60E-05	2.62E-06	2.00E-05
<b>AC073218.1</b>	rs892457	2:34667721	intergenic	3.43E-08	3.92E-05	3.89E-05	3.92E-05	3.89E-05	3.56E-06	2.71E-05
-	rs2357442	14:52607967	intergenic	4.40E-08	5.03E-05	4.98E-05	5.03E-05	4.98E-05	4.19E-06	3.19E-05
<b>PC1 (EV=1391.99)</b>										
<b>AGXT2</b>	<b>rs37369</b>	<b>5:35037115</b>	<b>missense</b>	<b>2.49E-11</b>	<b>2.85E-08</b>	<b>2.85E-08</b>	<b>2.85E-08</b>	<b>2.85E-08</b>	<b>2.85E-08</b>	<b>2.17E-07</b>
<b>FOXN3</b>	<b>rs390345</b>	<b>14:89976534</b>	<b>intronic</b>	<b>7.46E-11</b>	<b>8.55E-08</b>	<b>8.54E-08</b>	<b>8.55E-08</b>	<b>8.54E-08</b>	<b>4.27E-08</b>	<b>3.26E-07</b>
FAM49A	rs6741412	2:16815759	intronic	4.67E-10	5.35E-07	5.34E-07	5.35E-07	5.34E-07	1.43E-07	1.09E-06
between LOC100131241 and LOC124685	rs10512572	17:69512099	intergenic	4.99E-10	5.71E-07	5.70E-07	5.71E-07	5.70E-07	1.43E-07	1.09E-06
EYA1	rs79867447	8:72127562	intronic	7.46E-10	8.54E-07	8.51E-07	8.54E-07	8.51E-07	1.50E-07	1.14E-06
LMNA	rs12076700	1:156055099	intronic	7.87E-10	9.01E-07	8.97E-07	9.01E-07	8.97E-07	1.50E-07	1.14E-06
PCDH15	rs10825273	10:55968685	intronic	1.04E-09	1.19E-06	1.19E-06	1.19E-06	1.19E-06	1.70E-07	1.30E-06
upstream to BMP4	rs942316	14:54440983	intergenic	2.66E-09	3.04E-06	3.02E-06	3.04E-06	3.02E-06	3.80E-07	2.90E-06
TEX41	rs10496971	2:145769943	intronic	3.71E-09	4.25E-06	4.22E-06	4.25E-06	4.22E-06	4.72E-07	3.60E-06
between HAS2-AS1 and MRPS36P3	rs7844723	8:122908503	intergenic	8.11E-09	9.29E-06	9.21E-06	9.29E-06	9.21E-06	9.29E-07	7.08E-06
EYA2	rs58733120	20:45803852	intronic	2.19E-08	2.51E-05	2.49E-05	2.51E-05	2.49E-05	2.14E-06	1.63E-05
FAM49A	rs11096686	2:16815892	intronic	2.25E-08	2.57E-05	2.55E-05	2.57E-05	2.55E-05	2.14E-06	1.63E-05
EYA1	rs73684719	8:72131359	intronic	2.62E-08	3.00E-05	2.96E-05	3.00E-05	2.96E-05	2.30E-06	1.76E-05
XXYLT1	rs950257	3:194847650	intronic	3.94E-08	4.51E-05	4.46E-05	4.51E-05	4.46E-05	3.22E-06	2.46E-05

Based on assumption that polygenic traits are likely to be dominated by numerous alleles with small causal effect we also report the GWAS-level associations as these may represent true associations. Significantly associated markers that meet the p-value threshold of 5.00E-08 are located in genes or near genes responsible for various cellular functions. In general, these factors can be arbitrarily divided into three main categories: 1) genes with known roles in the craniofacial morphogenesis and/or mutated in various hereditary syndromes displaying craniofacial abnormalities; 2) genes or pseudo-genes without known function in the craniofacial morphology regulation or previously uncharacterized genes; and 3) non-protein coding genes, such as lncRNA class genes. There are also a number of significant variants that are located in the intergenic regions, with or without proximity to open reading frames (ORFs).

The majority of associated markers (n=32) are located in 23 protein-coding genes and pseudo-genes such as AGXT2, ASTN2, BCOR, BMP4, CACNB4, COL11A1, DCT, EGFR, EYA1, EYA2, FAM49A, FGF14, FOXN3, KIAA1217, LHX8, LMNA, MBD5, MEIS2, MYO5A, PAX3, PCDH15, RTTN, SMAD1, STON1, TMTC2, XXYLT1 and ZEB1.

Four variants are present in RNA-coding (lncRNA) genes, which include AC073218.1, RP11-494M8.4, RP11-408B11.2, and TEX41.

The rest of the markers (n=10) are found in the intergenic regions, near the following genes and pseudogenes: CTAGE13P, FRMD1, HAS2-AS1, LOC124685, LOC100131241, MRPS36P3, PTCH1, SLC25A5P2, SV2B and TRNAY16P.

Analysis of the functional annotation of significant markers revealed that one SNP represent missense mutation, one SNP is a synonymous transversion, 31 markers are located in intronic sequences and 13 markers are located in intergenic regions (Table 2). The majority of significantly associated SNPs (n=30) are found in the regulatory elements of the genome, such as in transcription factor (TF) binding sites, and represent potentially functional SNPs (pfSNPs). These variants may be involved in “fine tuning” of the normal craniofacial phenotype as part of the enhancer/silencer mechanisms, as has been recently suggested [77].

The nasal area measurements, using either “n”, “prn”, “sn” or “al” landmarks, produced the majority of the total number of significant associations (9 out of 11). These measurements include nasal width (al-al), nasal tip protrusion (sn-prn), nasion depth angle (zy\_l -n-zy\_r), nasolabial angle (prn-sn-ls), two similar measurements of the transverse nasal prominence angles (t\_l-prn-t\_r) and (zy\_l-prn-zy\_r), nasal index (al-al/n-sn), nasal tip protrusion – width

index (sn-prn/al-al), and nose face width index (al-al/zy-zy). The apparent overrepresentation of associations with the nasal area may be a result of the easier location and consequent superior reproducibility of the nasal area landmark measurements on 3D images. It may also be the result of specific selection of candidate genes from the JAX mice database resource, which focused on mutants that displayed various nasal area abnormalities. Notably, the two transverse nasal prominence angles provide similar anthropometric information. These essentially duplicate measurements were collected in order to compensate for losing data because of reduced quality images. In fact, both measurements showed significant associations with 8 shared markers, while all nasal area measurements in total share 18 significantly associated markers. (Table 2).

The analysis of direct cranial measurements and their relative indices revealed significant associations only of the Cephalic index (CI) with 9 SNPs.

The association analysis of the principal components (PC) representing all the craniofacial measurements, revealed one principal component that was associated with 14 genetic markers (Table 2). Additional three SNPs demonstrated some borderline associations at the p-values between  $6.15E-08$  and  $9.52E-08$  (Table 2). These variants include rs7460457 (COL22A1); rs79867447 (EYA1); rs1481800 (EYA1), which were associated with the ex-ex, ls-sto and eu-eu distances respectively. While we are not making strong claims for the significance of these variants, we believe this information would be useful for subsequent replication studies.

## **Craniofacial gene and SNP annotations**

The following section summarizes the genetic association results, providing brief annotation of the significantly associated genes and SNPs. Functional annotations, such as predicted molecular function, link to a biological process and a protein class of the 23 protein-coding genes and pseudo-genes (AGXT2, ASTN2, BCOR, BMP4, CACNB4, COL11A1, DCT, EGFR, EYA1, EYA2, FAM49A, FGF14, FOXN3, FMN1, KIAA1217, LMNA, MBD5, MEIS2, MYO5A, PAX3, PCDH15, RTTN, SMAD1, STON1, XXYL1 and ZEB1) have been visualised using the PANTHER resource [78] and summarized in supplemental materials (Supplemental Figs. S29-S31).



## Significantly associated genes with previously demonstrated role in craniofacial morphogenesis and/or mutated in hereditary syndromes displaying craniofacial abnormalities

A potentially functional SNP rs2289266 in the intron of the Paired Box 3 gene (**PAX3**) was associated with the transverse nasal prominence angle 2 (p-value 1.95E-08). This gene is a member of the paired box (PAX) family of transcription factors, which play critical roles during foetal development. The PAX3 protein regulates cell proliferation, migration and apoptosis. Mutations in PAX3 are associated with Waardenburg syndrome (OMIM: 193500), which is characterized by a prominent and broad nasal root, a round or square nose tip, hypoplastic alae, increased lower facial height and other craniofacial abnormalities.

Notably, three other SNPs in this gene, rs974448, rs7559271 and rs1978860, were previously associated with normal variability of the nasion position [22] and the distance between the eyeballs and the nasion [20]. None of these SNPs was included in this study, as a result of primer design failure. No LD between rs2289266 and any of the previously associated markers in the PAX3 gene was detected. Nevertheless, the association of another variant in the PAX3 gene can be considered an independent confirmation of this gene's involvement in regulation of normal craniofacial morphology.

SNPs rs4908280 and rs11164649 which are located in the regulatory element of the Collagen gene (**COL11A1**) intronic sequence, were associated with the cephalic index (p-values 1.66E-09 and 1.70E-08 respectively). COL11A1 encodes one of the two alpha chains of type XI fibrillar collagen and is known to have multiple transcripts as a result of alternative splicing. The secreted protein is hypothesised to play an important role in fibrillogenesis by controlling lateral growth of collagen II fibrils. Another potentially functional SNP, rs7460457, which is located in the intron of the **COL22A1** gene from the same Collagen family, was associated with the interexocanthal width distance (at the borderline p-value of 9.52E-08).

Notably, the same polymorphism (rs11164649) was recently linked to normal-range effects in various craniofacial traits, specifically eyes, orbits, nose tip, lips, philtrum and lateral parts of the mandible, although the measurements of the cephalic index were not performed in this study [93]. Our findings should be considered as independent confirmation of COL11A2

gene and its specific polymorphism rs11164649 involvement in shaping the normal craniofacial morphology.

According to the MGI database, transgenic mice with shortened COL11A2 mRNA (the second alpha chain of type XI fibrillar collagen) display abnormal facial phenotypes, including a triangular face and shorter and dimpled nasal bones [30]. Interestingly, COL11A1 and other Collagen family genes were found to be mutated in Stickler (OMIM: 604841) and Marshall Syndromes (OMIM: 154780). These two inherited disorders display very similar phenotypes and each is characterized by a distinctive facial appearance, with flat midface, very small jaw, cleft lip/palate, large eyes, short upturned nose, eye abnormalities, round face and short stature. However, the facial features of Stickler syndrome are less severe and include a flat face with depressed nasal bridge and cheekbones, caused by underdeveloped bones in the middle of the face. Another member of the collagen family, COL17A1, was recently associated with the distance between the eyeballs and the nasion [23]. Our finding of genetic associations of additional members of the Collagen family provides further evidence of the importance of polymorphisms in these genes in determining the normal variety of specific craniofacial features.

Intergenic SNP rs942316, which is located upstream to the Bone Morphogenetic Protein (**BMP4**) gene, was strongly associated with the following traits: nasion depth angle (p-value 3.74E-14), transverse nasal prominence 1 angle (p-value 5.51E-11) and PC1 (p-value 2.66E-09). The BMP4 gene is a transforming growth factor, belonging to the beta superfamily, which includes large families of growth and differentiation factors. This gene plays an important role in the onset of endochondral bone formation in humans, including induction of cartilage and bone formation and specifically tooth development and limb formation. Gene ontology annotations related to this gene include heparin binding and cytokine activity. BMP4 mutations have been associated with a variety of bone diseases, including orofacial cleft 11 (OMIM: 600625), Fibrodysplasia Ossificans (OMIM: 135100) and microphthalmia syndromic 6 (OMIM: 607932).

rs2476230 was associated with the transverse nasal prominence 1 angle (p-value 3.10E-08). This SNP is located in the intron of the **FGF14** gene, which encodes the Fibroblast Growth Factor 14 protein. FGF family members possess broad mitogenic and cell survival activities and are involved in a variety of biological processes, such as embryonic development, cell growth, morphogenesis and tissue repair. Mutations in various genes from the FGF family are associated with Ladd syndrome (OMIM:149730), the symptoms of which include various

craniofacial dysmorphisms such as high forehead, cleft lip and palate and cup-shaped, small ears.

SNP rs2290332 represents a synonymous variant in the Myosin VA (Heavy Chain 12, Myosin) gene (**MYO5A**). This variant was associated with cephalic index (p-value 5.56E-10) and nasion depth angle (p-value 3.67E-10).

**MYO5A** is one of three myosin V heavy-chain genes, belonging to the myosin gene superfamily. Myosin V is a class of actin-based motor proteins involved in cytoplasmic vesicle transport and anchorage, spindle-pole alignment and mRNA translocation. It mediates the transport of vesicles to the plasma membrane, including melanosome transport. Mutations in this gene were associated with a number of neuroectodermal diseases, such as Griscelli syndrome. Additional mutations in this gene were associated with a rare inherited condition Piebaldism (OMIM:172800). The symptoms of Piebaldism include partial albinism and anomalies of the mouth area development, such as lips and philtrum abnormalities. Despite being a “silent” mutation, rs2290332 is located in the POLR2A TF binding site and may therefore affect various processes such as transcription, translation, splicing and mRNA transport, as has been shown in other studies [79].

SNP rs7569399 was associated with the nasion depth angle (p-value 2.39E-10). This variant is located in the intron of the Methyl-CpG Binding Domain Protein 5 (**MBD5**). Encoded protein binds to the heterochromatin and potentially regulates cell division, growth and differentiation. Mutations in this gene were previously associated with different hereditary syndromes, including microcephaly.

Variant rs12041465, which is located in the intron of LIM Homeobox 8 (**LHX8**) was associated with transverse nasal prominences 1 and 2 angles (p-values 2.46E-10 and 2.30E-08 respectively). **LHX8** is a transcription factor and a member of the LIM homeobox family of proteins, which are involved in patterning and differentiation of various tissue types. Mutations in this gene were associated with clefts of the secondary palate in mouse model [80, 81].

Intronic SNP rs5963728 was associated with the nasion depth angle (p-value 3.67E-10) and transverse nasal prominence 1 angle (p-value 1.28E-08). This marker is located in the BCL6 Corepressor (**BCOR**). This gene encodes a zinc finger transcription repressor, which may specifically inhibit gene expression leading to apoptosis influence when recruited to promoter regions by sequence-specific DNA-binding proteins such as BCL6 and MLLT3.

Interestingly, mutations in **BCOR** were associated with oculofaciocardiodental syndrome (OMIM:300166). The symptoms of this syndrome include various facial anomalies such as long narrow face, dip-set small or absent eyes, broad nasal tip, which is also divided by cleft and missing or abnormal teeth.

Three intronic SNPs in the Eyes Absent Homolog 1 (**EYA1**) gene were associated with several craniofacial traits. The variant rs79867447 was associated with the nose width (p-value 3.92E-08), nasal index (p-value 3.53E-08), nose-face width index (p-value 1.10E-09), PC1 (p-value 7.46E-10) and upper vermilion height (borderline p-value 6.15E-08). The variant rs1481800 was associated with the cephalic index (p-value 2.07E-09) and maximum head width (borderline p-value 8.85E-08). The variant rs73684719 was found in association with PC1 (p-value 2.62E-08). All three variants are located in the potentially regulatory elements of the genome and are likely to affect TF binding sites. No linkage disequilibrium has been detected between these markers.

The EYA1 encoded protein functions as histone phosphatase, regulating transcription during organogenesis in kidney and various craniofacial features such as branchial arches, eye and ear. EYE1 mutated mice display various craniofacial anomalies of the inner ear, mandible, maxilla and reduced skull [30]. Mutations in the human ortholog have been associated with several craniofacial conditions such as otofaciocervical syndrome (OMIM:166780), Weyers acrofacial dysostosis (OMIM:193530) and branchiootic syndrome (OMIM:608389).

Intronic SNP rs58733120 was associated with the nose width (p-value 5.37E-10), nasal index (p-value ), nasal tip protrusion – nose width index (p-value 9.46E-09), nose-face width index (p-value 3.38E-09) and PC1 (p-value 2.19E-08) phenotypes. This variant is located in the regulatory element of the **EYA2** gene, which belongs to the same eyes absent protein family as EYA1. The EVC2 gene encodes a positive regulator of the hedgehog signalling pathway and plays a critical role in bone formation and skeletal development. The EVC2 gene was named after the Ellis Van Creveld Syndrome 2 (OMIM:225500), which is associated with this gene [82]. This disorder is characterised by various skeletal abnormalities and particularly by very short stature.

SNP rs12076700 in the intron of the Lamin A gene (**LMNA**) was associated with the nasal tip protrusion – nose width index (p-value 1.66E-09), nasion depth angle (p-value 1.09E-09), transverse nasal prominences 1 and 2 angles (p-values 1.30E-09 and 1.54E-09) and PC1 (p-value 7.87E-10).

**LMNA**, together with other Lamin proteins, is a component of a fibrous layer on the nucleoplasmic side of the inner nuclear membrane, which provides a framework for the nuclear envelope and also interacts with chromatin. LMNA encoded protein acts to disrupt mitosis and induces DNA damage in vascular smooth muscle cells, leading to mitotic failure, genomic instability, and premature senescence of the cell. This gene has been found mutated in Mandibuloacral Dysplasia which is characterized by various skeletal and craniofacial abnormalities, including delayed closure of the cranial sutures and undersized jaw [83].

Marker rs74884233 was associated with the transverse nasal prominences 1 and 2 angles (p-values 1.40E-09 and 1.20E-08 respectively). This variant is located in the intron of the Rotatin gene (**RTTN**). RTTN gene is involved in the maintenance of normal ciliary structure, which in turn effects the developmental process of left-right organ specification, axial rotation, and perhaps notochord development.

SNP rs17020235 was associated with the nasolabial angle (p-value 2.07E-09). This potentially functional variant is located in the intron of the SMAD Family Member 1 gene (**SMAD1**). SMAD1 is a transcriptional modulator activated by BMP (bone morphogenetic proteins) type 1 receptor kinase, which is involved in a range of biological activities including cell growth, apoptosis, morphogenesis, development and immune responses.

SMAD1 mutant mice display anterior truncation of the head with only one brachial arch present. In human, SMAD1 mutations (together with RUNX2), are associated with the Cleidocranial Dysplasia (OMIM:119600), which is a Craniosynostosis-type disorder affecting cranial bones, palate and other tissues.

Potentially functional SNP rs3910620 was associated with the nasion depth angle (p-value 2.16E-09) and transverse nasal prominence 1 (p-value 1.06E-08). This variant is located in intron of the **KIAA1217** gene. This gene is likely an ortholog of the mouse Enhancer Trap Locus 4, which is required for normal development of intervertebral disks. Mutations in the KIAA1217 gene were associated with a musculo-skeletal disorder called Lumbar disc herniation (OMIM:603932).

Variant rs2122497 in the intron of the Meis Homeobox 2 gene (**MEIS2**) was associated with the nasal tip protrusion - width index (p-value 2.33E-09). This gene encodes a homeobox protein belonging to the TALE ('three amino acid loop extension') family of homeodomain-containing proteins. TALE homeobox proteins are highly conserved transcription regulators, shown to regulate embryonic developmental. MEIS2 gene is believed to be involved in the

transcriptional activation of the ELA1 and EPHA8 enhancers in the developing midbrain. Multiple transcript variants encoding distinct isoforms have been described for this gene. MEIS2 was associated with the 15q14 Microdeletion syndrome with symptoms that include asymmetric skull with narrow forehead, bulbous nasal tip, high nasal bridge and eye abnormalities [84].

SNP rs57585041 was associated with the nose width (p-value 6.05E-09). This variant is located in the TF binding site downstream to Patched 1 (**PTCH1**) gene. The encoded protein is the receptor for sonic hedgehog (SHH), indian hedgehog (IHH) and desert hedgehog (DHH), which are the key factors in embryonic morphogenesis. Mice homozygous for PTCH1 have multiple skeletal defects, including a shortened face, wide set eyes, excessive skin and thickened foot pads.

SNP rs950257 was associated with the PC1 trait (p-value 3.94E-08). This intronic variant is located in the **XXYL1** gene, which codes for Xyloside Xylosyltransferase 1. This protein is an Alpha-1,3-xylosyltransferase, which elongates the O-linked xylose-glucose disaccharide attached to EGF-like repeats in the extracellular domain of Notch proteins signalling network. Notch proteins are the key regulators of embryonic development, which demonstrate a highly conserved sequence in various species. Interestingly, mutations in Notch proteins are associated with Hajdu–Cheney syndrome (OMIM:10250) and Alagille syndrome (OMIM:118450). The main phenotypic symptoms of these conditions include various malformations of the craniofacial tissues, including broad, prominent forehead, deep-set eyes and a small pointed chin.

SNP rs17335905 was associated with the nose-face width index (p-value 4.74E-08). This potentially functional variant is located in the intron of the **EGFR** gene, which encodes the Epidermal Growth Factor Receptor. EGFR is a cell surface protein that binds to epidermal growth factor (EGF). Binding of the protein to a ligand induces activation of several signalling cascades and leads to cell proliferation, cytoskeletal rearrangement and anti-apoptosis. Mouse carrying mutations in EGFR, express short mandible and cleft palate.



## Significantly associated SNPs, located in genes or pseudo-genes that were not linked to craniofacial morphology regulation or genes with unknown function

Intronic variant rs59037879 in the Zinc Finger E-Box Binding Homeobox 1 (**ZEB1**) was found associated with cephalic index (p-value 6.27E-10), nasion depth angle (p-value 5.77E-11) and transverse nasal prominences angles 1 and 2 (p-values 5.65E-09 and 5.31E-12 respectively). This gene encodes a zinc finger transcription factor, which is a transcriptional repressor. It regulates expression of different genes, such as interleukin-2 (IL-2) gene, ATPase transporting polypeptide (ATP1A1) gene and E-cadherin (CDH1) promoter in various cell types and also represses stemness-inhibiting microRNA. Mutations in this gene were previously associated with Corneal Dystrophy and various types of cancer.

A missense mutation rs37369 in the Alanine--Glyoxylate Aminotransferase 2 gene (**AGXT2**) was associated with nose width (p-value 1.04E-09), nasal index: (p-value 1.25E-10), nasal tip protrusion – nose width index: (p-value 1.09E-11), nasion depth angle (p-value 1.31E-10), transverse nasal prominence 1 and 2 angles (p-values 8.65E-10 and 1.46E-09 respectively) and PC1 (p-value 2.49E-11). This protein plays an important role in regulating blood pressure in the kidney through metabolizing asymmetric dimethylarginine (ADMA), which is an inhibitor of nitric-oxide (NO) synthase.

Two intronic SNPs located in the regulatory element of the Calcium Channel Voltage-Dependent Beta 4 Subunit (**CACNB4**) gene intron, were significantly associated with numerous traits. rs16830498 was associated with the cephalic index (p-value 7.57E-11) and transverse nasal prominence 1 angle (p-value 6.16E-10)

The beta subunit of voltage-dependent calcium channels may increase peak calcium current by shifting the voltage dependencies of activation and inactivation, modulating G protein inhibition and controlling the alpha-1 subunit membrane targeting. CACNB4 may be expressed in different isoforms through alternative splicing. Certain mutations in this gene have been associated with various forms of epilepsy, although no association with normal or abnormal craniofacial variation has been previously reported.

Potentially functional intronic SNPrs10825273 located in the regulatory elements of the Protocadherin-Related 15 (**PCDH15**) gene, was found in association with several facial phenotypes, such as the nasion depth angle (p-value 1.55E-12), nasal tip protrusion – width



index (p-value 2.82E-09), cephalic index (p-value 9.93E-09), transverse nasal prominence 1 (p-value 3.36E-10) and PC1 (p-value 1.04E-09). PCDH15 is a member of the cadherin superfamily, which encodes an integral membrane protein that mediates calcium-dependent cell-cell adhesion and is known to have numerous alternative splicing variants. It plays an essential role in the maintenance of normal retinal and cochlear function. Mutations in this gene result in hearing loss and are associated with Usher Syndrome Type IIA (OMIM: 276901).

Two intronic variants in the Family With Sequence Similarity 49 Member A gene (**FAM49A**) were associated with multiple craniofacial traits. rs6741412 was found in association with the nasion depth angle (p-value 3.00E-11), transverse nasal prominence 1 and 2 angles (p-values 3.22E-12 and 2.75E-09 respectively) and PC1 (p-value 4.67E-10). rs11096686 was associated with the nasal tip protrusion - width index (p-value 1.72E-09) and PC1 (p-value 2.25E-08). The FAM49A protein is known to interact with hundreds of miRNA molecules during pre-implantation of the mouse embryo and also expressed in the developing chick wing, but no information on its specific function or disease association have been identified.

SNP rs390345, located in the intronic regulatory sequence of the Forkhead Box N3 gene (**FOXN3**), was associated with the nasal tip protrusion – nose width index: (p-value 1.96E-10) and PC1 (p-value 7.46E-11). FOXN3 produces multiple splicing variants and acts as a transcriptional repressor. It is proposed to be involved in DNA damage-inducible cell cycle arrests at G1 and G2. There are no previous reports on FOXN3 association with either normal craniofacial development or pathological conditions.

SNP rs1407995 variant, located in the regulatory region of the Dopachrome Tautomerase gene (**DCT**) was associated with the nasion depth angle (p-value 9.42E-11) and transverse nasal prominence 1 angle (p-value 6.49E-10). Another pfSNP rs1325611 in the same gene was found in association with the same traits: nasion depth angle (p-value 3.85E-10) and transverse nasal prominence 1 angle (p-value 2.12E-09). DCT is known to be involved in regulating eumelanin and pheomelanin levels. According to GO annotation, this gene might be involved in the oxidoreductase and dopachrome isomerase activities. Spontaneous mutation in the DCT locus cause limb deformity in mice.

Genetic variant rs7574549 marker was associated with the nasion depth angle (p-value 1.33E-09) and transverse nasal prominence 1 angle (p-value 1.27E-10). This variant is

located at the 3' prime UTR sequence of the Stonin 1 gene (**STON1**). This sequence is a part of the naturally occurring read-through products of the neighbouring STON1 and TFIIA-alpha and beta-like factor (GTF2A1L) genes. STON1 gene encodes one of two human homologs of the *Drosophila melanogaster* stoned B protein, which may be involved in the endocytic machinery of the plasma membrane. GTF2A1L is a germ cell-specific transcription factor that is able to stabilize the binding of TBP to DNA. Co-transcription of this gene and the neighbouring upstream STON1 generates a rare transcript, which encodes a fusion protein comprised of sequence sharing identity with each individual gene product. No craniofacial abnormalities were associated with either STON1 or GTF2A1L genes in human or animal models.

SNP rs10513300 was associated only with the nasion depth angle (p-value  $2.37 \times 10^{-10}$ ) distance. This intronic variant is located in the Astrotactin 2 (**ASTN2**), which encodes a protein that is expressed in the brain and may regulate neuronal migration. A deletion at this locus has been associated with schizophrenia.

SNP rs11115552 was associated with the transverse nasal prominence angle 1 (p-value  $6.10 \times 10^{-8}$ ). This intronic variant is located in the regulatory sequence of the **TMTC2** gene, which encodes a Transmembrane and Tetratricopeptide Repeat Containing 2 protein, which might be required for Calcium ion homeostasis.

### **Significantly associated SNPs located in the non-protein coding genes, such as lncRNA class genes**

Intronic SNP rs10496971 in the **TEX41** (Testis Expressed 41) gene demonstrated the lowest p-value ( $1.86 \times 10^{-17}$ ) among the craniofacial trait associations. This marker was associated with the nasion depth angle, a measurement that in general produced the most significant associations. The same SNP produced significant associations with nasal tip protrusion – width index (p-value  $9.36 \times 10^{-11}$ ), transverse nasal prominence angles 1 and 2 (p-values  $6.80 \times 10^{-14}$  and  $5.52 \times 10^{-9}$  respectively), cephalic index (p-value  $5.315 \times 10^{-9}$ ) and PC1 (p-value  $3.71 \times 10^{-9}$ ).

TEX41 is a long intergenic non-protein coding RNA (lncRNA) class gene, which is located on chromosome 2 and has 43 transcript variants as a result of alternative splicing. lncRNAs are known as regulators of diverse cellular processes. However, the function of this gene remains unknown. Despite its name, this gene is expressed in a variety of tissues, with the highest demonstrated levels in kidney. Its potential involvement in craniofacial genetics, and specifically in influencing normal facial variation, has not been reported previously. Notably, the rs10496971 variant is located in the regulatory element of the genome (as well as 49 other associated SNPs) and may influence normal craniofacial morphology by affecting either enhancer or silencer sequences or transcriptional factor (TF) binding sites [77].

SNP rs1871428, located in the regulatory sequence of the lncRNA class gene **RP11-503C24.4**, whose specific function is unknown, was associated with the nasion depth angle (p-value 2.70E-10) and transverse nasal prominence 1 angle (p-value 1.37E-08).

The SNP rs1482795, located in the RNA gene **RP11-494M8.4**, was associated with the nose width (p-value 7.68E-10) and nasal index (p-value 1.83E-08).

Two SNPs rs892457 and rs892458, located in the non-protein coding lncRNA gene **AC073218.1**, were associated with the transverse nasal prominence 2 angle (p-value 3.43E-08) and (p-value 1.73E-08), respectively.

SNP rs7311798, located in the lncRNA gene **RP11-408B11.2** was associated with the nose width (p-value 1.66E-07).

SNP rs7844723 in the **RP11-785H20.1** (lncRNA gene) was associated with the nasion depth angle (p-value 1.56E-09) and PC1 (p-value 8.11E-09) phenotypes.

SNP rs2357442 was associated with the transverse nasal prominence 2 angle (p-value 4.40E-08). This variant is located in the Long Interspersed Nuclear Element 1 (**LINE-1**) retrotransposon sequence, which in turn shows homology with uncategorized mRNA KC832805 on the Y-chromosome.

LINE-1 elements comprise approximately 21% of the reference genome, and have been shown to modulate expression and produce novel splice isoforms of transcripts from genes that span or neighbour the LINE-1 insertion site. In addition, rs2357442 is located close to three pseudo-genes with unknown function: SLC25A5P2, LOC100130842 and RP11-1033H12.1, while the last two represent RNA-coding lncRNA genes.

## Significantly associated SNPs located in the intergenic regions

SNP rs10512572, located between Serpine1 MRNA Binding Protein 1 pseudogene (**LOC100131241**) and Myosin, Light Chain 6 Alkali Smooth Muscle and Non-Muscle pseudogene (**LOC124685**), was associated with nasal tip protrusion (p-value 2.22E-08), nasal tip protrusion – nose width index (p-value 4.48E-11), nasion depth angle (p-value 2.89E-12), transverse nasal prominence 1 angle (p-value 1.50E-11), transverse nasal prominence 2 angle (p-value 1.38E-11) and PC1 (p-value 4.99E-10). While pseudogenes in general are non-protein coding, their sequences can be functional and play important roles in different biological processes [85]. It should be noted that some genes may be incorrectly defined as pseudogenes, based solely on their sequence computational analysis [86]. The function of these two pseudogene sequences is unknown.

SNP rs8035124 was significantly associated with the nose width (p-value 1.52E-10). This variant is located between the Synaptic Vesicle Glycoprotein 2B (**SV2B**) and Transfer RNA Tyrosine 16 (Anticodon GUA) Pseudogene (**TRNAY16P**) genes. The SV2B is a protein coding gene, which plays a role in the control of regulated secretion in neural and endocrine cells. The TRNAY16P is a pseudogene with unknown function.

Additional SNP rs373272 was associated with cephalic index (p-value 2.40E-08) . However, no genes were identified within 50 kb window of its chromosomal location.

## Genetic associations of non-anthropometric traits

The eyelid is a non-anthropometric, although visible facial trait. The formation of eyelids in the foetus starts at approximately 8 weeks when the folds of surface ectoderm overgrow the eyes to form the eyelids, which remain closed until the seventh month of development. The eyelid can be largely segregated into single or double (according to presence or absence of a skin fold), with the former more dominant in the East Asian populations and the latter in the remaining world populations.

Analysis of the HWE-filtered data did not produce any significant associations with the eyelid phenotype. Analysis of the HWE non-filtered data however, produced associations with two polymorphisms. Two variants, rs4823810 in the CELSR1 gene and rs11217807 in the POU2F3 gene, demonstrated genetic associations at unadjusted p-values of 4.21E-08 and 2.46E-07 respectively. The CELSR1 gene belongs to the cadherin superfamily, involved in

cell adhesion and receptor-ligand interactions. In general, cadherins participate in signal transduction in the WNT signalling pathway, which is central to embryonic development. This specific protein is a developmentally regulated neural-specific factor, which plays an unspecified role in early embryogenesis. The AmiGO ontology database suggests however, that this gene may play a role in the establishment of planar cell polarity. Specifically, it may be one of several factors that coordinate organization of groups of cells in the plane of an epithelium, such that they all orient in a similar direction. It may be hypothesized that it could play a role in regulating facial symmetry. Interestingly, the CELSR1 deficiency in mice results in various craniofacial defects, such as craniorachischisis (neural tube defects) and failure of the eyelid closure [87, 88].

The POU2F3 gene encodes a transcriptional factor, which regulates cell type-specific differentiation pathways. The encoded protein is primarily expressed in the epidermis and plays a critical role in keratinocyte proliferation and differentiation, while also being a candidate tumour suppressor protein. The association of POU2F3 with the eyelid or any other craniofacial trait has not been reported before.

The association analysis of the ear lobe phenotype (attached/detached) did not produce associations above the significance threshold level. While the association analysis of the craniofacial traits has been focused on the HWE-filtered data, the association results of the non-filtered data cannot be neglected given the association results of pigmentation traits. Despite possible population stratification, these results may indicate true associations, rather than false positives. A replication study using larger sample size of various population groups may indeed clarify this question.

## Conclusions

This study focused on the identification of genetic markers in a set of candidate genes associated with various craniofacial traits, representing the most comprehensive scan for genetic markers involved in normal craniofacial development performed to date. We identified 11 phenotypes that were significantly associated (unadjusted p-value  $\leq 5.00\text{E-}08$ ) with 42 genomic variants in 26 genes and 15 intergenic regions. Following the application of over-conservative Bonferroni correction (p-value threshold of  $1.6\text{E-}07$ ), associations were observed between 8 craniofacial traits and 12 SNPs located in 12 genes and intergenic regions. We reported all the significant markers that met the less stringent GWAS threshold,

as complex trait such as craniofacial morphology is likely to be influenced by a large number of alleles with relatively small individual effect, similar to height [89, 90].

The association of the PAX3 gene and two genes from the collagen family (COL11A1 and COL22A1) with several craniofacial traits confirm previous findings [11, 22, 23, 91]. In fact, an intronic SNP rs11164649 that was associated with cephalic index in the current study was recently associated with normal-range effects in various craniofacial traits and used for their prediction [91]. However, the majority of genetic associations are novel. These include 27 significantly associated markers in protein-coding genes and pseudo-genes, , such as AGXT2, ASTN2, BCOR, BMP4, CACNB4, DCT, EGFR, EYA1, STON1, TMTC2 and ZEB1. Additionally, 19 SNPs in intergenic regions adjacent to several genes, such as BMP4,DMD, FAM47A, FRMD1, LOC124685, LOC100131241, SV2Bwere found to be associated with craniofacial morphology. Some of these genes were previously linked to craniofacial embryogenesis, while others represent novel factors.

Seven further genetic variants were found in lncRNA genes, which have not been linked to craniofacial morphogenesis before. Notably, rs10496971 in the TEX41 lncRNA gene demonstrated the strongest association among craniofacial trait associations, including the nasion depth angle (p-value 1.86E-17). These findings suggest that there may be a yet unexplored level of epigenetic regulation affecting craniofacial morphology. lncRNAs are a recently discovered class of factors, whose expression is thought to be important for the regulation of gene expression through several different mechanisms involving competition with transcription by recruitment of specific epigenetic factors to promoter regions, as well as indirectly affecting gene expression by interacting with miRNA and other cellular factors [92]. The comprehensive role of epigenetic regulation in general, and in craniofacial embryonic development in particular, is poorly understood. There is a limited number of recent studies revealing thousands of enhancer sequences, predicted to be active in the developing craniofacial complex in mice [77, 93] and potentially in humans. Both the epistatic and epigenetic interactions may represent a more complex level of craniofacial morphology regulation and require further investigation.

Even though a relatively high number of phenotypes were studied (92 linear and angular measurements and indices), this may still represent an oversimplification of the complexity of the human face. Despite the importance of the association between specific 3D measurements and SNPs demonstrated in this study, the association of facial shapes, represented by the

principal components should better represent the face. Given that embryonic developmental processes such as cell proliferation, polarity orientation and migration occur in a 3D environment, principal components that in essence denote specific facial shapes, may provide a more accurate representation of these processes. However, only one of the 20 principle components showed significant associations at the GWAS threshold level. While the explanation of this observation is unclear, it is consistent with other similar studies [22, 23]. The specific anthropometric measurements on the other hand, produced numerous significant associations, identifying many genes and intergenic regions that appear to play important roles in the development of normal human facial appearance. While replication of these results is critical, it was not performed yet due to time and budget limitations of this study. Nevertheless, the findings from this preliminary study add significantly to unravelling the genetic basis of craniofacial morphology and the data are freely available for a further analysis upon request.

Given the high complexity of the face, as well as the composite nature of the genetic regulation that affects its development, alternative comprehensive approaches of capturing facial morphology would be beneficial. A number of such methods has recently revealed additional genes with specific polymorphisms associated with the development of craniofacial traits within the normal variation range [91, 94]. Further studies may involve the use of these or alternative methods to capture the majority of variation in craniofacial traits. Craniofacial phenotypes, together with additional external visible traits such as sex, age and BMI and ancestry, could be treated as a “vector”, which could then be used to predict appearance [95].

A recent attempt to predict facial appearance was performed using only 24 SNPs [96]. This approach has promise, although it is largely based on reconstruction of a ‘facial composite image’ through prediction of ancestry, sex, pigmentation and human perception of faces. This approach is reasonable, but it does not negate the use of association studies looking at specific craniofacial traits. Genetic association studies of a large scope of individual anthropometric measurements are essential to provide information on specific genes and their polymorphisms, which affect these traits and may therefore be useful in predicting the size and the shape of specific facial features.

Additional association studies on large sample sizes, incorporating dense SNP panels or whole genome sequencing approaches, in conjunction with either a comprehensive set of anthropometrical measurements or morphologically adequate representation of the



craniofacial characteristics would be a valuable adjunct to the promising results obtained in this study. These studies will not only improve our understanding of the genetic factors regulating craniofacial morphology, but will also enable a better prediction of the visual appearance of a person from DNA.

## List of abbreviations

3D: 3-Dimensional; AIMS: ancestry informative markers; ASW: African ancestry in Southwest USA; BAM: Binary alignment map; BMI: Body Mass Index; CAU: Caucasian; CHB: Han Chinese in Beijing, China; ENIGMA: ENtire Genome INterface for Exploring SnpS; EVT: Externally visible characteristic; DVI: Disaster victim identification; FDP: Forensic DNA phenotyping; GWAS: Genome wide association studies; HWE: Hardy-Weinberg equilibrium; JPT: Japanese in Tokyo, Japan; LD: linkage disequilibrium; lncRNAs: long non-coding RNAs; LINE-1: Long Interspersed Nuclear Element 1; MAF: Minor allele frequency; measurement error; ME: Measurement error; MD: Mean difference; OMIM: Online Mendelian Inheritance in Man; ORFs: open reading frames; pfsNP: Potentially functional SNP; PCA: Principal component analysis; RGB: Red, Green, Blue (colours); SNP: Single-nucleotide polymorphism; SNAP: SNP Annotation and Proxy Search; STR: Short tandem repeat; TF: Transcription factor; VCF: Variant Call Format; YRI: Yoruba in Ibadan, Nigeria.

## Declarations

## Ethics (and consent to participate)

The participants provided their written informed consent to participate in this study, which was approved by the Bond University Ethics committee (RO-510).

## Competing interests

The authors declare that they have no competing interests.



# Authors' contributions

MB designed the study, carried out the molecular genetic studies, carried out the data analysis, participated in the statistical analysis and drafted the manuscript. PB performed the statistical analysis and drafted the manuscript. AvD participated in the design of the study and drafted the manuscript. All authors read and approved the final manuscript.

# Consent to Publish

Not applicable

# Availability of data and materials

The genomic data supporting the conclusions of this article are included within the article and its additional files.

# Funding

The funding for this research was provided by the Technical Support Working Group (Award Number: IS-FB-2946) and Pelerman Holdings Pt Ltd.

# Acknowledgments

We would like to thank the volunteers who participated in this study without whom we could not have performed this research. We would like to thank Olga Kondrashova who helped with the sample collection. We also thank Technical Support Working Group (TSWG) and Pelerman Holdings Pt Ltd for their generous support of this project.

# References

1. Kohn LAP. The Role of Genetics in Craniofacial Morphology and Growth. Annual Review of Anthropology. 1991;20:261-78. doi: 10.2307/2155802.
2. Sperber GH, Sperber SM, Guttman GD. Craniofacial embryogenetics and development. 2nd ed. Shelton, CT: People's Medical Pub. House USA; 2010. 250p. ill. (some colour) p.
3. Richtsmeier JT, Cheverud JM. Finite element scaling analysis of human craniofacial growth. Journal of craniofacial genetics and developmental biology. 1986;6(3):289-323. PubMed PMID: 3771738.
4. Neubauer S, Gunz P, Hublin JJ. The pattern of endocranial ontogenetic shape changes in humans. J Anat. 2009;215(3):240-55. doi: 10.1111/j.1469-7580.2009.01106.x. PubMed PMID: 19531085; PubMed Central PMCID: PMC2750758.
5. Sturm RA. Molecular genetics of human pigmentation diversity. Human Molecular Genetics. 2009;18(R1):R9-R17. doi: 10.1093/hmg/ddp003.
6. Shkoukani MA, Chen M, Vong A. Cleft Lip - A Comprehensive Review. Front Pediatr. 2013;1:53. doi: 10.3389/fped.2013.00053. PubMed PMID: 24400297; PubMed Central PMCID: PMC3873527.
7. Kimonis V, Gold J-A, Hoffman TL, Panchal J, Boyadjiev SA. Genetics of Craniosynostosis. Semin Pediatr Neurol. 2007;14:150-61.
8. Weinberg SM, Naidoo SD, Bardi KM, Brandon CA, Neiswanger K, Resick JM, et al. Face shape of unaffected parents with cleft affected offspring: combining three-dimensional surface imaging and geometric morphometrics. Orthodontics & Craniofacial Research. 2009;12(4):271-81. doi: 10.1111/j.1601-6343.2009.01462.x.
9. Anna K Coussens CRW, Ian P Hughes, C, Phillip Morris, Angela van Daal, Peter J Anderson, Barry C Powell. Unravelling the molecular control of calvarial suture fusion in children with craniosynostosis. BMC Genomics. 2007;8:458.
10. Coussens AK, Hughes IP, Wilkinson CR, Morris CP, Anderson PJ, Powell BC, et al. Identification of genes differentially expressed by prematurely fused human sutures using a novel in vivo[thin space]-[thin space]in vitro approach. Differentiation. 2008;76(5):531-45. doi: DOI: 10.1111/j.1432-0436.2007.00244.x.
11. Boehringer S, van der Lijn F, Liu F, Gunther M, Sinigerova S, Nowak S, et al. Genetic determination of human facial morphology: links between cleft-lips and normal

variation. *Eur J Hum Genet.* 2011;19(11):1192-7. doi: 10.1038/ejhg.2011.110. PubMed PMID: 21694738; PubMed Central PMCID: PMC3198142.

12. Nakamura A, Hattori M, Sakaki Y. A novel gene isolated from human placenta located in Down syndrome critical region on chromosome 21. *DNA research : an international journal for rapid publication of reports on genes and genomes.* 1997;4(5):321-4. doi: 10.1093/dnares/4.5.321. PubMed PMID: 9455479.

13. El Ghouzzi V. Mutations in the basic domain and the loop-helix II junction of TWIST abolish DNA binding in Saethre-Chotzen syndrome. *FEBS Lett.* 2001;492:112-8.

14. Kamath BM, Stolle C, Bason L, Colliton RP, Piccoli DA, Spinner NB, et al. Craniosynostosis in Alagille syndrome. *Am J Med Genet.* 2002;112(2):176-80. doi: 10.1002/ajmg.10608. PubMed PMID: 12244552.

15. Hood RL, Lines MA, Nikkel SM, Schwartzentruber J, Beaulieu C, Nowaczyk MJ, et al. Mutations in SRCAP, encoding SNF2-related CREBBP activator protein, cause Floating-Harbor syndrome. *Am J Hum Genet.* 2012;90(2):308-13. doi: 10.1016/j.ajhg.2011.12.001. PubMed PMID: 22265015; PubMed Central PMCID: PMC3276662.

16. Roper RJ, Baxter LL, Saran NG, Klinedinst DK, Beachy PA, Reeves RH. Defective cerebellar response to mitogenic Hedgehog signaling in Down's syndrome mice. *Proc Natl Acad Sci U S A.* 2006;103(5):1452-6.

17. Croonen EA, van der Burgt I, Kapusta L, Draaisma JM. Electrocardiography in Noonan syndrome PTPN11 gene mutation--phenotype characterization. *American journal of medical genetics Part A.* 2008;146A(3):350-3. doi: 10.1002/ajmg.a.32140. PubMed PMID: 18203203.

18. Fourie Z, Damstra J, Gerrits PO, Ren Y. Evaluation of anthropometric accuracy and reliability using different three-dimensional scanning systems. *Forensic Sci Int.* 2011;207(1-3):127-34. doi: 10.1016/j.forsciint.2010.09.018. PubMed PMID: 20951517.

19. Toma AM, Zhurov A, Playle R, Ong E, Richmond S. Reproducibility of facial soft tissue landmarks on 3D laser-scanned facial images. *Orthodontics & Craniofacial Research.* 2009;12(1):33-42. doi: 10.1111/j.1601-6343.2008.01435.x.

20. Kovacs L, Zimmermann A, Brockmann G, Baurecht H, Schwenzer-Zimmerer K, Papadopoulos NA, et al. Accuracy and Precision of the Three-Dimensional Assessment of the Facial Surface Using a 3-D Laser Scanner. *IEEE Transactions on Medical Imaging.* 2006;25(6):742-54. PubMed PMID: 21197761.

21. Coussens AK, Daal Av. Linkage disequilibrium analysis identifies an FGFR1 haplotype-tag SNP associated with normal variation in craniofacial shape. *Genomics*. 2005;85(5):563-73. doi: DOI: 10.1016/j.ygeno.2005.02.002.
22. Paternoster L, Zhurov AI, Toma AM, Kemp JP, St Pourcain B, Timpson NJ, et al. Genome-wide association study of three-dimensional facial morphology identifies a variant in PAX3 associated with nasion position. *Am J Hum Genet*. 2012;90(3):478-85. Epub 2012/02/22. doi: 10.1016/j.ajhg.2011.12.021. PubMed PMID: 22341974; PubMed Central PMCID: PMC3309180.
23. Liu F, van der Lijn F, Schurmann C, Zhu G, Chakravarty MM, Hysi PG, et al. A genome-wide association study identifies five loci influencing facial morphology in Europeans. *PLoS Genet*. 2012;8(9):e1002932. doi: 10.1371/journal.pgen.1002932. PubMed PMID: 23028347; PubMed Central PMCID: PMC3441666.
24. Michel S, Liang L, Depner M, Klopp N, Ruether A, Kumar A, et al. Unifying candidate gene and GWAS Approaches in Asthma. *PLoS One*. 2010;5(11):e13894. doi: 10.1371/journal.pone.0013894. PubMed PMID: 21103062; PubMed Central PMCID: PMC2980484.
25. Sun J, Jia P, Fanous AH, Webb BT, van den Oord EJ, Chen X, et al. A multi-dimensional evidence-based candidate gene prioritization approach for complex diseases-schizophrenia as a case. *Bioinformatics*. 2009;25(19):2595-6602. doi: 10.1093/bioinformatics/btp428. PubMed PMID: 19602527; PubMed Central PMCID: PMC2752609.
26. Hrdy DB. Analysis of hair samples of mummies from Semma South (Sudanese Nubia). *Am J Phys Anthropol*. 1978;49(2):277-82. doi: 10.1002/ajpa.1330490217. PubMed PMID: 717558.
27. Fitzpatrick TB. The validity and practicality of sun-reactive skin types I through VI. *Archives of Dermatology*. 1988;124(6):869.
28. Sturm RA, Larsson M. Genetics of human iris colour and patterns. *Pigment Cell Melanoma Res*. 2009;22(5):544-62. doi: 10.1111/j.1755-148X.2009.00606.x. PubMed PMID: 19619260.
29. Farkas LG. *Anthropometry of the head and face*. 2nd ed. New York: Raven Press; 1994. xix, 405 p. p.
30. Site MMRW. The Jackson Laboratory, Bar Harbor, Maine. World Wide Web (<http://mousemutant.jax.org/>) Bar Harbor, Maine.2010 [cited 2013 March]. Available from: <http://mousemutant.jax.org/>.

31. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet.* 2000;25(1):25-9. doi: 10.1038/75556. PubMed PMID: 10802651; PubMed Central PMCID: PMC3037419.
32. Online Mendelian Inheritance in Man O. McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University (Baltimore, MD) [cited 2012 March]. Available from: <http://omim.org/>.
33. Rebhan M, ChalifaCaspi V, Prilusky J, Lancet D. GeneCards: Integrating information about genes, proteins and diseases. *Trends in Genetics.* 1997;13(4):163-. doi: Doi 10.1016/S0168-9525(97)01103-7. PubMed PMID: WOS:A1997WQ96900011.
34. Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, et al. Genome-wide detection and characterization of positive selection in human populations. *Nature.* 2007;449(7164):913-8. doi: 10.1038/nature06250. PubMed PMID: 17943131; PubMed Central PMCID: PMC2687721.
35. Zhou N, Wang L. Effective selection of informative SNPs and classification on the HapMap genotype data. *BMC Bioinformatics.* 2007;8(1):484. PubMed PMID: doi:10.1186/1471-2105-8-484.
36. Kosoy R, Nassir R, Tian C, White PA, Butler LM, Silva G, et al. Ancestry informative marker sets for determining continental origin and admixture proportions in common populations in America. *Hum Mutat.* 2009;30(1):69-78. doi: 10.1002/humu.20822. PubMed PMID: 18683858; PubMed Central PMCID: PMC3073397.
37. Paschou P, Lewis J, Javed A, Drineas P. Ancestry informative markers for fine-scale individual assignment to worldwide populations. *Journal of Medical Genetics.* doi: 10.1136/jmg.2010.078212.
38. Londin ER, Keller MA, Maista C, Smith G, Mamounas LA, Zhang R, et al. CoAIMs: a cost-effective panel of ancestry informative markers for determining continental origins. *PLoS One.* 2010;5(10):e13443. doi: 10.1371/journal.pone.0013443. PubMed PMID: 20976178; PubMed Central PMCID: PMC2955551.
39. Bulbul O, Filoglu G, Altuncul H, Aradas AF, Ruiz Y, Fondevila M, et al. A SNP multiplex for the simultaneous prediction of biogeographic ancestry and pigmentation type. *Forensic Sci Inter: Genet Suppl.* 2011;3(1):e500-e1. doi: 10.1016/j.fsigss.2011.10.001.
40. Tandon A, Patterson N, Reich D. Ancestry informative marker panels for African Americans based on subsets of commercially available SNP arrays. *Genet Epidemiol.* 2011;35(1):80-3. doi: 10.1002/gepi.20550. PubMed PMID: 21181899.

41. Phillips C, Freire Aradas A, Kriegel AK, Fondevila M, Bulbul O, Santos C, et al. Eurasiaplex: a forensic SNP assay for differentiating European and South Asian ancestries. *Forensic Sci Int Genet.* 2013;7(3):359-66. doi: 10.1016/j.fsigen.2013.02.010. PubMed PMID: 23537756.
42. Gettings KB, Lai R, Johnson JL, Peck MA, Hart JA, Gordish-Dressman H, et al. A 50-SNP assay for biogeographic ancestry and phenotype prediction in the U.S. population. *Forensic Sci Int Genet.* 2014;8(1):101-8. doi: 10.1016/j.fsigen.2013.07.010. PubMed PMID: 24315596.
43. Kidd KK, Speed WC, Pakstis AJ, Furtado MR, Fang R, Madbouly A, et al. Progress toward an efficient panel of SNPs for ancestry inference. *Forensic Sci Int Genet.* 2014;10C(0):23-32. doi: 10.1016/j.fsigen.2014.01.002. PubMed PMID: 24508742.
44. Amigo J, Salas A, Phillips C. ENGINES: exploring single nucleotide variation in entire human genomes. *BMC Bioinformatics.* 12(1):105. PubMed PMID: doi:10.1186/1471-2105-12-105.
45. Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A method and server for predicting damaging missense mutations. *Nat Methods.* 2010;7(4):248-9. doi: 10.1038/nmeth0410-248. PubMed PMID: 20354512; PubMed Central PMCID: PMC2855889.
46. Marinescu V, Kohane I, Riva A. MAPPER: a search engine for the computational identification of putative transcription factor binding sites in multiple genomes. *BMC Bioinformatics.* 2005;6(1):79. PubMed PMID: doi:10.1186/1471-2105-6-79.
47. Wang J, Ronaghi M, Chong SS, Lee CGL. pfSNP: An integrated potentially functional SNP resource that facilitates hypotheses generation through knowledge syntheses. *Human Mutation.* 32(1):19-24. doi: 10.1002/humu.21331.
48. Coassin S, Brandstätter A, Kronenberg F. Lost in the space of bioinformatic tools: A constantly updated survival guide for genetic epidemiology. *The GenEpi Toolbox. Atherosclerosis.* 209(2):321-35. doi: DOI: 10.1016/j.atherosclerosis.2009.10.026.
49. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype data. *Genetics.* 2000;155(2):945-59. PubMed PMID: WOS:000087475100039.
50. Thorisson G, Smith A, Krishnan L, Stein L. The International HapMap Project Web site. *Genome Res.* 2005;15:1592.
51. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum*

Genet. 2007;81(3):559-75. doi: 10.1086/519795. PubMed PMID: 17701901; PubMed Central PMCID: PMC1950838.

52. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet.* 2006;38(8):904-9. doi: 10.1038/ng1847. PubMed PMID: 16862161.

53. Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet.* 2006;2(12):e190. doi: 10.1371/journal.pgen.0020190. PubMed PMID: 17194218; PubMed Central PMCID: PMC1713260.

54. UniProt C. Activities at the Universal Protein Resource (UniProt). *Nucleic Acids Res.* 2014;42(Database issue):D191-8. doi: 10.1093/nar/gkt1140. PubMed PMID: 24253303; PubMed Central PMCID: PMC3965022.

55. Rappaport N, Nativ N, Stelzer G, Twik M, Guan-Golan Y, Stein TI, et al. MalaCards: an integrated compendium for diseases and their annotation. *Database : the journal of biological databases and curation.* 2013;2013:bat018. doi: 10.1093/database/bat018. PubMed PMID: 23584832; PubMed Central PMCID: PMC3625956.

56. Warde-Farley D, Donaldson SL, Comes O, Zuberi K, Badrawi R, Chao P, et al. The GeneMANIA prediction server: biological network integration for gene prioritization and predicting gene function. *Nucleic Acids Res.* 2010;38(Web Server issue):W214-20. doi: 10.1093/nar/gkq537. PubMed PMID: 20576703; PubMed Central PMCID: PMC2896186.

57. Reference Genome Group of the Gene Ontology C. The Gene Ontology's Reference Genome Project: a unified framework for functional annotation across species. *PLoS Comput Biol.* 2009;5(7):e1000431. doi: 10.1371/journal.pcbi.1000431. PubMed PMID: 19578431; PubMed Central PMCID: PMC2699109.

58. Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, et al. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Research.* 2001;29(1):308-11. doi: 10.1093/nar/29.1.308.

59. Durbin RM, Abecasis GR, Altshuler DL, Auton A, Brooks LD, Gibbs RA, et al. A map of human genome variation from population-scale sequencing. *Nature.* 467(7319):1061 - 73. PubMed PMID: doi:10.1038/nature09534.

60. Chelala C, Khan A, Lemoine NR. SNPnexus: a web database for functional annotation of newly discovered and public domain single nucleotide polymorphisms. *Bioinformatics.* 2009;25(5):655-61. doi: 10.1093/bioinformatics/btn653. PubMed PMID: 19098027; PubMed Central PMCID: PMC2647830.



61. Rajeevan H, Cheung KH, Gadagkar R, Stein S, Soundararajan U, Kidd JR, et al. ALFRED: an allele frequency database for microevolutionary studies. *Evol Bioinform Online*. 2005;1:1-10. PubMed PMID: 19325849; PubMed Central PMCID: PMC2658869.
62. Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ, de Bakker PI. SNAP: a web-based tool for identification and annotation of proxy SNPs using HapMap. *Bioinformatics*. 2008;24(24):2938-9. doi: 10.1093/bioinformatics/btn564. PubMed PMID: 18974171; PubMed Central PMCID: PMC2720775.
63. Boyle AP, Hong EL, Hariharan M, Cheng Y, Schaub MA, Kasowski M, et al. Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res*. 2012;22(9):1790-7. doi: 10.1101/gr.137323.112. PubMed PMID: 22955989; PubMed Central PMCID: PMC3431494.
64. Aung SC, Ngim RC, Lee ST. Evaluation of the laser scanner as a surface measuring tool and its accuracy compared with direct facial anthropometric measurements. *British journal of plastic surgery*. 1995;48(8):551-8. PubMed PMID: 8548155.
65. Bianchi SD, Spada MC, Bianchi L, Verz   L, Vezzetti E, Tornincasa S, et al. Evaluation of scanning parameters for a surface colour laser scanner. *International Congress Series*. 2004;1268:1162-7. doi: 10.1016/j.ics.2004.03.264. PubMed PMID: 13327565.
66. Kusnoto B, Evans CA. Reliability of a 3D surface laser scanner for orthodontic applications. *American journal of orthodontics and dentofacial orthopedics : official publication of the American Association of Orthodontists, its constituent societies, and the American Board of Orthodontics*. 2002;122(4):342-8. PubMed PMID: 12411877.
67. Ma L, Xu T, Lin J. Validation of a three-dimensional facial scanning system based on structured light techniques. *Computer Methods & Programs in Biomedicine*. 2009;94(3):290-8. doi: 10.1016/j.cmpb.2009.01.010. PubMed PMID: 37572488.
68. Gwilliam JR, Cunningham SJ, Hutton T. Reproducibility of soft tissue landmarks on three-dimensional facial scans. *European journal of orthodontics*. 2006;28(5):408-15. doi: 10.1093/ejo/cjl024. PubMed PMID: 16901962.
69. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, et al. GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol*. 2010;28(5):495-501. doi: 10.1038/nbt.1630. PubMed PMID: 20436461.
70. Weston EM, Friday AE, Lio P. Biometric evidence that sexual selection has shaped the hominin face. *PLoS One*. 2007;2(8):e710. doi: 10.1371/journal.pone.0000710. PubMed PMID: 17684556; PubMed Central PMCID: PMC1937021.



71. Cooper RS, Tayo B, Zhu X. Genome-wide association studies: implications for multiethnic samples. *Hum Mol Genet.* 2008;17(R2):R151-5. Epub 2008/10/15. doi: 10.1093/hmg/ddn263. PubMed PMID: 18852204; PubMed Central PMCID: PMC2782359.
72. Barnholtz-Sloan JS, Chakraborty R, Sellers TA, Schwartz AG. Examining population stratification via individual ancestry estimates versus self-reported race. *Cancer Epidemiol Biomarkers Prev.* 2005;14(6):1545-51. doi: 10.1158/1055-9965.EPI-04-0832. PubMed PMID: 15941970.
73. Tang H, Quertermous T, Rodriguez B, Kardia SL, Zhu X, Brown A, et al. Genetic structure, self-identified race/ethnicity, and confounding in case-control association studies. *Am J Hum Genet.* 2005;76(2):268-75. doi: 10.1086/427888. PubMed PMID: 15625622; PubMed Central PMCID: PMC1196372.
74. Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population stratification in genome-wide association studies. *Nat Rev Genet.* 2010;11(7):459-63. Epub 2010/06/16. doi: 10.1038/nrg2813. PubMed PMID: 20548291; PubMed Central PMCID: PMC2975875.
75. Panagiotou OA, Ioannidis JP, Genome-Wide Significance P. What should the genome-wide significance threshold be? Empirical replication of borderline genetic associations. *International journal of epidemiology.* 2012;41(1):273-86. doi: 10.1093/ije/dyr178. PubMed PMID: 22253303.
76. Dudbridge F, Gusnanto A. Estimation of significance thresholds for genomewide association scans. *Genet Epidemiol.* 2008;32(3):227-34. doi: 10.1002/gepi.20297. PubMed PMID: 18300295; PubMed Central PMCID: PMC2573032.
77. Attanasio C, Nord AS, Zhu Y, Blow MJ, Li Z, Liberton DK, et al. Fine tuning of craniofacial morphology by distant-acting enhancers. *Science.* 2013;342(6157):1241006. doi: 10.1126/science.1241006. PubMed PMID: 24159046.
78. Thomas PD, Campbell MJ, Kejariwal A, Mi H, Karlak B, Daverman R, et al. PANTHER: a library of protein families and subfamilies indexed by function. *Genome Res.* 2003;13(9):2129-41. doi: 10.1101/gr.772403. PubMed PMID: 12952881; PubMed Central PMCID: PMC403709.
79. Goymer P. Synonymous mutations break their silence. *Nat Rev Genet.* 2007;8(2):92-.
80. Zhao Y, Guo YJ, Tomac AC, Taylor NR, Grinberg A, Lee EJ, et al. Isolated cleft palate in mice with a targeted mutation of the LIM homeobox gene *lhx8*. *Proc Natl Acad Sci U S A.* 1999;96(26):15002-6. PubMed PMID: 10611327; PubMed Central PMCID: PMC24762.

81. Zhang Y, Mori T, Takaki H, Takeuch M, Iseki K, Hagino S, et al. Comparison of the expression patterns of two LIM-homeodomain genes, Lhx6 and L3/Lhx8, in the developing palate. *Orthod Craniofac Res*. 2002;5(2):65-70. PubMed PMID: 12086327.
82. Galdzicka M, Patnala S, Hirshman MG, Cai JF, Nitowsky H, A Egeland J, et al. A new gene, EVC2, is mutated in Ellis-van Creveld syndrome. *Molecular genetics and metabolism*. 2002;77(4):291-5. doi: 10.1016/s1096-7192(02)00178-6.
83. Novelli G, Muchir A, Sangiuolo F, Helbling-Leclerc A, D'Apice MR, Massart C, et al. Mandibuloacral dysplasia is caused by a mutation in LMNA-encoding lamin A/C. *Am J Hum Genet*. 2002;71(2):426-31. doi: 10.1086/341908. PubMed PMID: 12075506; PubMed Central PMCID: PMC379176.
84. van Bon BWM, Mefford HC, Menten B, Koolen DA, Sharp AJ, Nillesen WM, et al. Further delineation of the 15q13 microdeletion and duplication syndromes: a clinical spectrum varying from non-pathogenic to a severe outcome. *Journal of Medical Genetics*. 2009;46(8):511-23. doi: 10.1136/jmg.2008.063412.
85. Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. A coding-independent function of gene and pseudogene mRNAs regulates tumour biology. *Nature*. 2010;465(7301):1033-8.
86. Hirotsume S, Yoshida N, Chen A, Garrett L, Sugiyama F, Takahashi S, et al. An expressed pseudogene regulates the messenger-RNA stability of its homologous coding gene. *Nature*. 2003;423(6935):91-6.
87. Allache R, De Marco P, Merello E, Capra V, Kibar Z. Role of the planar cell polarity gene CELSR1 in neural tube defects and caudal agenesis. *Birth defects research Part A, Clinical and molecular teratology*. 2012;94(3):176-81. doi: 10.1002/bdra.23002. PubMed PMID: 22371354.
88. Meng Q, Jin C, Chen Y, Chen J, Medvedovic M, Xia Y. Expression of signaling components in embryonic eyelid epithelium. *PLoS One*. 2014;9(2):e87038. doi: 10.1371/journal.pone.0087038. PubMed PMID: 24498290; PubMed Central PMCID: PMC3911929.
89. Visscher PM. Sizing up human height variation. *Nature Genetics*. 2008;40(5):489-90. doi: 10.1038/ng0508-489. PubMed PMID: 18443579.
90. Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat Genet*. 2014;46(11):1173-86. doi: 10.1038/ng.3097. PubMed PMID: 25282103; PubMed Central PMCID: PMCPMC4250049.

91. Claes P, Liberton DK, Daniels K, Rosana KM, Quillen EE, Pearson LN, et al. Modeling 3D facial shape from DNA. *PLoS Genet.* 2014;10(3):e1004224. doi: 10.1371/journal.pgen.1004224. PubMed PMID: 24651127; PubMed Central PMCID: PMC3961191.
92. Lau E. Non-coding RNA: Zooming in on lncRNA functions. *Nat Rev Genet.* 2014;15(9):574-5. doi: 10.1038/nrg3795.
93. Wu H, Nord AS, Akiyama JA, Shoukry M, Afzal V, Rubin EM, et al. Tissue-Specific RNA Expression Marks Distant-Acting Developmental Enhancers. *PLoS Genet.* 2014;10(9):e1004610. doi: 10.1371/journal.pgen.1004610.
94. Peng S, Tan J, Hu S, Zhou H, Guo J, Jin L, et al. Detecting Genetic Association of Common Human Facial Morphological Variation Using High Density 3D Image Registration. *PLoS Comput Biol.* 2013;9(12):e1003375. doi: 10.1371/journal.pcbi.100337.
95. Wolf L, Donner Y, editors. An experimental study of employing visual appearance as a phenotype. *Computer Vision and Pattern Recognition, 2008 CVPR 2008 IEEE Conference on*; 2008: IEEE.
96. Claes P, Hill H, Shriver MD. Toward DNA-based facial composites: preliminary results and validation. *Forensic Sci Int Genet.* 2014;13:208-16. doi: 10.1016/j.fsigen.2014.08.008. PubMed PMID: 25194685.

# Figure legends, tables and additional file descriptions

**Figure 1. Anatomical position of the 32 manually annotated anthropometric landmarks used for calculation of linear and angular distances and ratios between the linear distances.** Some landmarks are not clearly visible due to image orientation. gn = Gnathion, pg= Pogonion, sl = Sublabiale, li = Labiale Inferius, sto = Stomion, ls = Labiale superius, ch-r = Chelion right, ch-l = Chelion left, go-r = Gonion Right, go-l = Gonion left, sn = Subnasale, prn= Pronasale, al-r = Alare right; al-l= Alare left, n = Nasion, g= Glabella; tr = Tragon, en-l = left Endocanthion, en-r = right Endocanthion, ex-r = Right Endocanthion; ex-l = left Endocanthion, ps-r = Palpebrale superius right, ps-l = Palpebrale superius left , pi-r = Palpebrale inferius right, pi-l = Palpebrale inferius left, zy-r = Zygion Right, zy-l = Zygion Left, pra-r = Tragon right, pra-l = Tragon Left, sba-l = Subalare left, sa-l = Superaurale Left, pa-l = Postaurale left.

**Figure 2. Illustration of linear and angular distances calculated from manually annotated landmark coordinates.**

**Figure 3. A simplified flow chart, summarizing major steps of the candidate gene search, SNP genotyping and data analysis workflow.**

**Figure 4. Population structure as represented by plotting genomic PCs 1 and 2, using 270 HapMap individuals.** YRI: Yoruba, Nigeria, Africa. JRI: Japanese, Tokyo, Japan. CHB: Han Chinese, Beijing, China. CEU: Utah residents with European ancestry.

## Additional files

**Figure S1.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of al-al distance and expected association based on the overall al-al distance distribution.

**Figure S2.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of ex-ex distance and expected association based on the overall ex-ex distance distribution.

**Figure S3.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of ls-sto distance and expected association based on the overall ls-sto distance distribution.

**Figure S4.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of sn-prn distance and expected association based on the overall sn-prn distance distribution.

**Figure S5.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of eu-eu distance and expected association based on the overall eu-eu distance distribution.

**Figure S6.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of cephalic index and expected association based on the overall cephalic index distribution.

**Figure S7.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of nasal index and expected association based on the overall nasal index distribution.

**Figure S8.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of nasal tip protrusion-width

index and expected association based on the overall nasal tip protrusion-width index distribution.

**Figure S9.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of nose-face width index and expected association based on the overall nose-face width index distribution.

**Figure S10.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of nasion depth angle and expected association based on the overall nasion depth angle distance distribution.

**Figure S11.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of nasolabial angle and expected association based on the overall nasolabial angle distance distribution.

**Figure S12.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of transverse nasal prominence angle 1 and expected association based on the overall nasal prominence angle 1 distribution.

**Figure S13.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of transverse nasal prominence angle 2 and expected association based on the overall transverse nasal prominence angle 2 distribution.

**Figure S14.** Q-Q plot of the PCA-corrected  $-\log_{10}$  p-values for the difference between the observed association for the tails of PC1 trait and expected association based on the overall PC1 trait distance distribution.

**Figure S15.** Manhattan plot of the genomic associations of the al-al distance, based on the initial p-values from analysis of the PCA-corrected data. The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue

line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S16. Manhattan plot of the genomic associations of the ex-ex distance, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S17. Manhattan plot of the genomic associations of the ls-sto distance, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S18. Manhattan plot of the genomic associations of the sn-prn distance, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S19. Manhattan plot of the genomic associations of the eu-eu distance, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue



line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S20. Manhattan plot of the genomic associations of the cephalic index, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S21. Manhattan plot of the genomic associations of the nasal index, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S22. Manhattan plot of the genomic associations of the nasal tip protrusion-width index, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S23. Manhattan plot of the genomic associations of the nose-face width index, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated

by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S24. Manhattan plot of the genomic associations of the nasion depth angle, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S25. Manhattan plot of the genomic associations of the nasolabial angle, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S26. Manhattan plot of the genomic associations of the transverse nasal prominence angle 1, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S27. Manhattan plot of the genomic associations of the transverse nasal prominence angle 2, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is

indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S28. Manhattan plot of the genomic associations of the PC1 trait, based on the initial p-values from analysis of the PCA-corrected data.** The  $-\log_{10}$  (P value) is plotted against the physical positions of each SNP on each chromosome. The basic significance threshold is indicated by the blue line for  $-\log_{10}(1e-5)$ , the candidate genes threshold for  $-\log_{10}(5e-7)$  is indicated by the green line and the genome-wide significance threshold for  $-\log_{10}(5e-8)$  is indicated by the red line.

**Figure S29. Pie chart, illustrating molecular function classification of the 37 human genes, harbouring 40 markers in association with craniofacial phenotypes.** The genes are: AGXT2, ASTN2, BCOR, BMP4, CACNB4, COL11A1, COL22A1, DCT, EGFR, EYA1, EYA2, FAM49A, FGF14, FOXN3, KIAA1217, LMNA, MBD5, MEIS2, MYO5A, PAX3, PCDH15, RTTN, SMAD1, STON1, XXYLT1 and ZEB1. Generated using PANTHER gene ontology resource.

**Figure S30. Pie chart, illustrating biological processes classification involving 23 human genes, harbouring 32 markers in association with craniofacial phenotypes.** The genes include: AGXT2, ASTN2, BCOR, BMP4, CACNB4, COL11A1, COL22A1, DCT, EGFR, EYA1, EYA2, FAM49A, FGF14, FOXN3, KIAA1217, LMNA, MBD5, MEIS2, MYO5A, PAX3, PCDH15, RTTN, SMAD1, STON1, XXYLT1 and ZEB1. Generated using PANTHER gene ontology resource.

**Figure S31. Pie chart, illustrating protein product classification of the 23 human genes, harbouring 32 markers in association with craniofacial phenotypes.** The genes are: AGXT2, ASTN2, BCOR, BMP4, CACNB4, COL11A1, COL22A1, DCT, EGFR, EYA1, EYA2, FAM49A, FGF14,

FOXN3, KIAA1217, LMNA, MBD5, MEIS2, MYO5A, PAX3, PCDH15, RTTN, SMAD1, STON1, XXYLT1 and ZEB1. Generated using PANTHER gene ontology resource.

**Table S1. Manually annotated facial landmarks used in the study.**

**Table S2. Genetic syndromes displaying various craniofacial abnormalities, used to locate candidate genes for the study.**

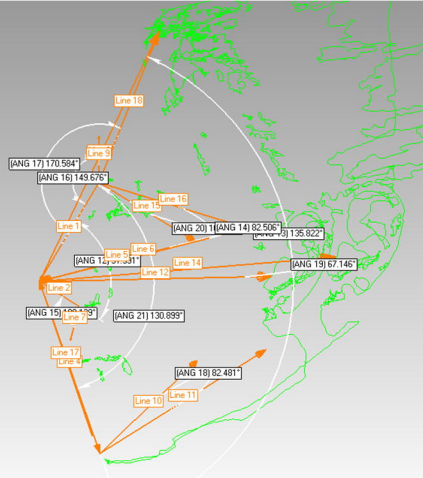
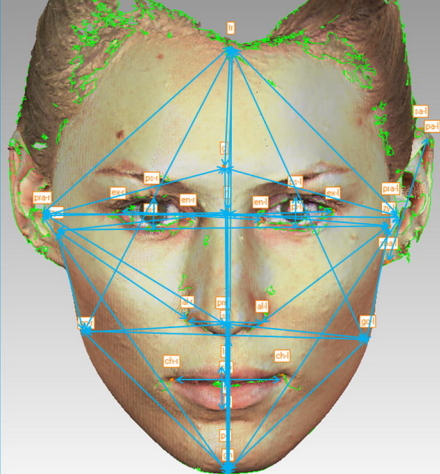
**Table S3. Genetic associations with pigmentation traits.** Gene: gene name; rs#: reference SNP ID number; SNP: chromosomal location of the marker; Genomic annotation: genomic location of the marker; UNADJ: Unadjusted p-values; BONF: Bonferroni single-step adjusted; HOLM: Holm (1979) step-down adjusted; SIDAK\_SS: Sidak single-step adjusted; SIDAK\_SD: Sidak step-down adjusted; FDR\_BH: Benjamini & Hochberg (1995) step-up FDR control; FDR\_BY: Benjamini & Yekutieli (2001) step-up FDR control.

**Appendix S1. A comprehensive list of web resources used for candidate gene search and its output.** Note the presence of multiple tabs in this spreadsheet.

**Appendix S2. Two spreadsheets, detailing a list of 8,518 genetic markers genotyped in 587 DNA samples and a list of 3,073 markers used for association analyses, following MAF (2%) filtering.**









# Candidate gene search, genotyping and genetic association results summary

496 Candidate genes and genomic regions identified through the literature and web resources search

Craniofacial disorders in human and model organisms, GWAS studies

Genes with SNPs showing high degree of differentiation between populations ( $F_{st}$ )

1,319 SNPs in 177 candidate genes and intergenic regions

High  $F_{st}$  SNPs

nsSNPs, TFB-sites SNPs, splicing sites, tag SNPs

+ 522 pigmentation markers

Anthropometric analysis of the 3-D facial images and physical measurements

6 cranial and 32 facial landmarks

92 linear and angular measurements and ratios, 10 principal components + 25 additional phenotypic traits and ancestry

587 DNA samples genotyped by MPS (Ion Torrent platform)

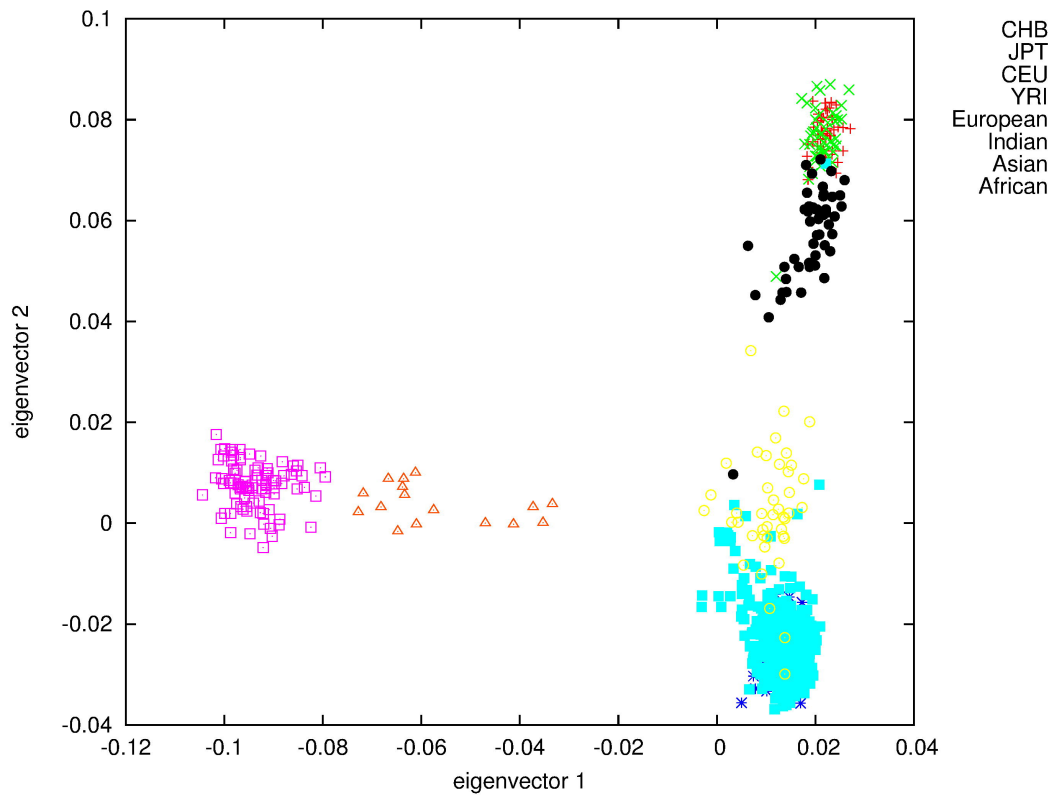
Custom Ampliseq assay

1,700 primer pairs, 32 barcodes per chip

6,945 SNPs ( $MAF \leq 1\%$ ) were genotyped

Genetic association analysis

Following QC filtering, 2,332 SNPs were analysed for potential association with the craniofacial traits and pigmentation. 19 craniofacial traits were strongly associated with 78 SNPs located in 36 genes and 14 intergenic regions at  $p\text{-value} \leq 5.00E-07$



# PANTHER Protein Class

Total # Genes: 37    Total # protein class hits: 40



[cell adhesion molecule \(P00009\)](#)

[c1 junction protein \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[cellular junction protein \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[cellular junction protein \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

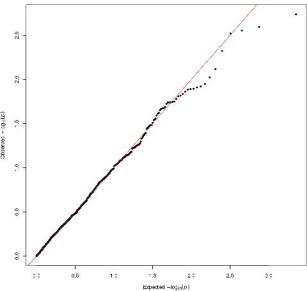
[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

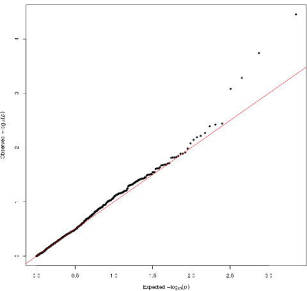
[caveolin-1 \(P00009\)](#)

[caveolin-1 \(P00009\)](#)

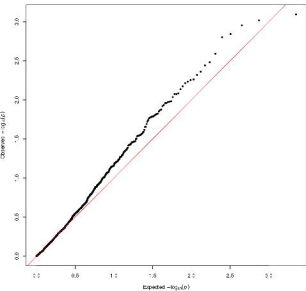
Noise width



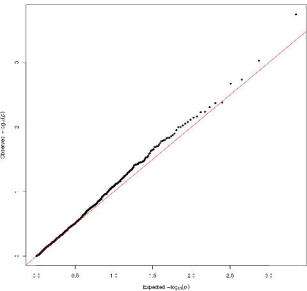
# Interecocal width



Upper vermilion height

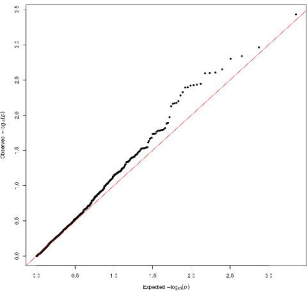


Nasal tip protrusion

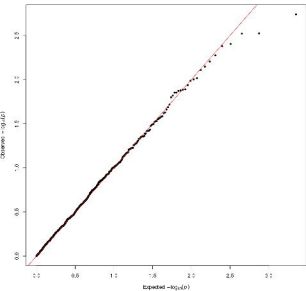




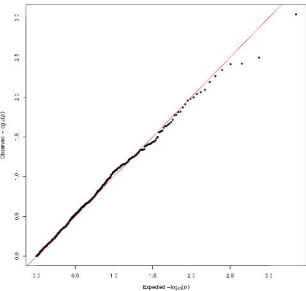
Maximum Head Width



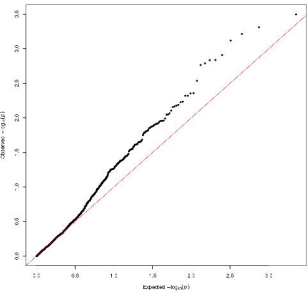
Cephalic index



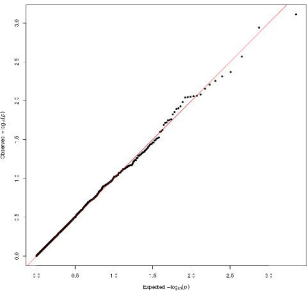
Nasal index



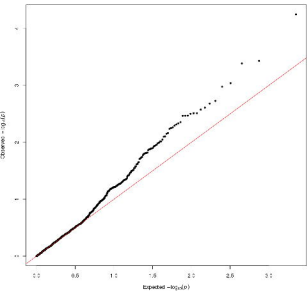
Nasal lip protrusion - nose width index



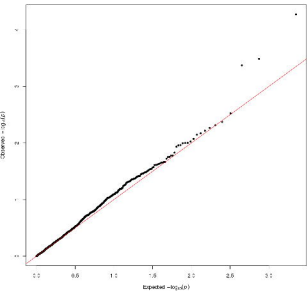
Nose-face width index



Naxon depth angle

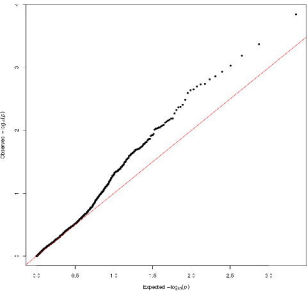


Noncollinear angle

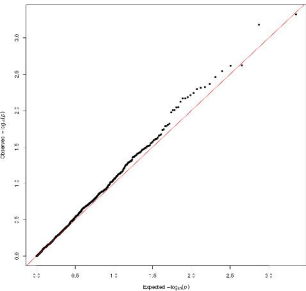




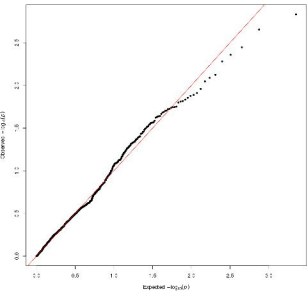
Transverse nasal prominence 1 angle



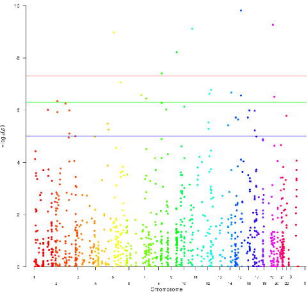
Transverse nasal prominence 2 angle



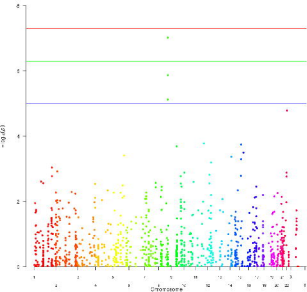
Grain EV = 2204.88



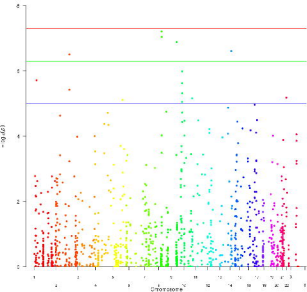
Noise width



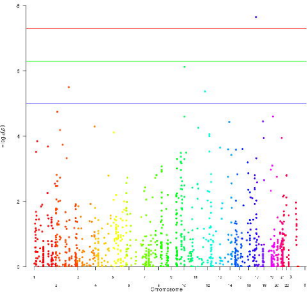
# Interocanthal width



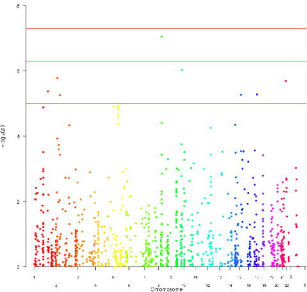
Upper vermilion height



Nasal tip protrusion

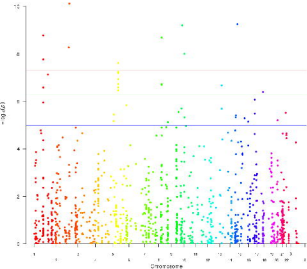


Maximum Head Width

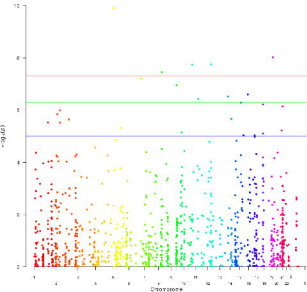




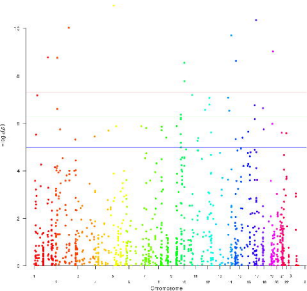
# Cephalic index



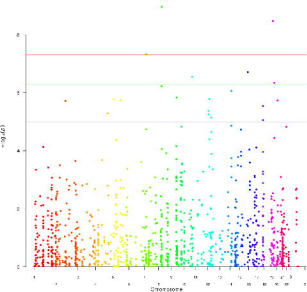
Nasal index



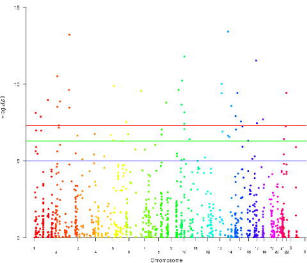
# Nasal tip protrusion - nose width index



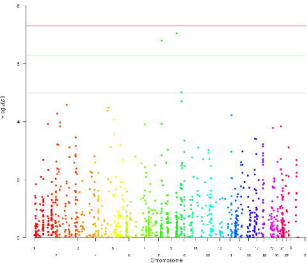
Rose-face width index



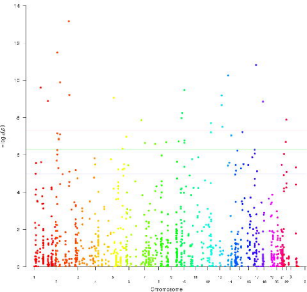
# Nasion depth angle



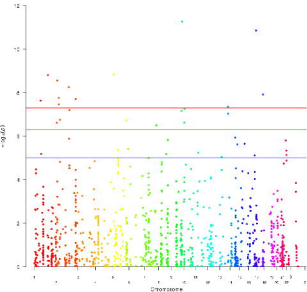
# Nucleolar angle



Transverse nasal prominence 1 angle

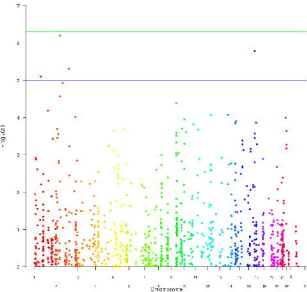


Transverse nasal prominence 2 angle





Cratio EV = 2294.88



## GO Molecular Functions

Total # Genes: 47    Total # Gene Ontology: 40



binding (GO:0005488)

catalytic activity (GO:0003674)

cytokine receptor activity (GO:0005115)

cytokine and chemokine receptor activity (GO:0005115)

receptor activity (GO:0005115)

structural molecular activity (GO:0005115)

transporter activity (GO:0005115)

GO Biological Processes  
Total # Genes: 37    Total # processes hit: 13



acoustic process (GO:0005945)

cellular response to stimulus (GO:0032502)

cellular response to stress (GO:0006954)

cellular component organization or biogenesis (GO:0071586)

cellular process (GO:0009957)

developmental process (GO:0032502)

response to stress (GO:0006954)

cellular response to hypoxia (GO:0001617)

cellular response to hypoxia (GO:0001617)

metabolic process (GO:0008152)

cellular response to hypoxia (GO:0001617)

response to stimulus (GO:0032502)